# The Next 25 Years of Computer Architecture?
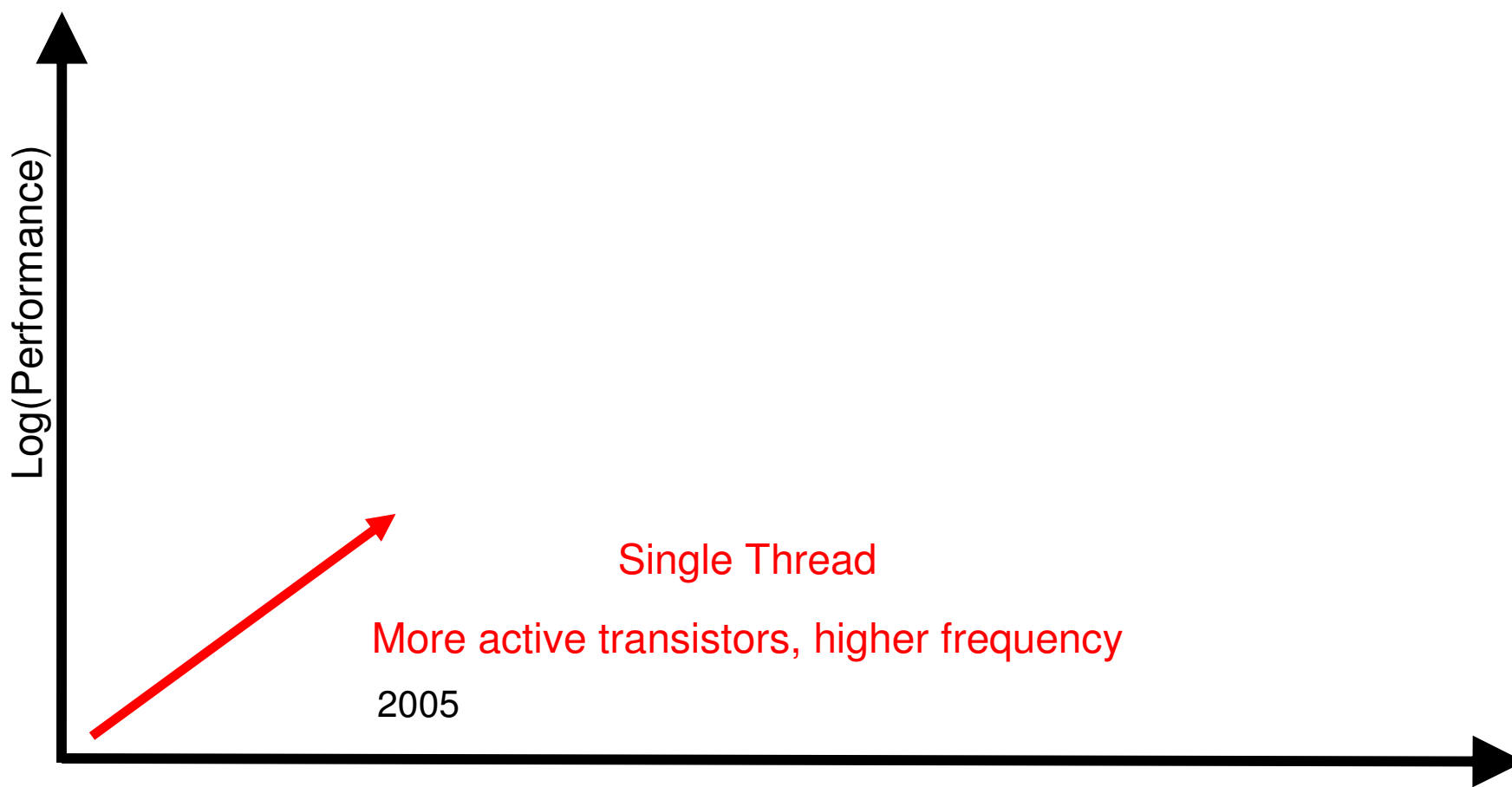
**HPPC 2009, Delft, August 25 2009**

## H. Peter Hofstee
## Cell/B.E. Chief Scientist
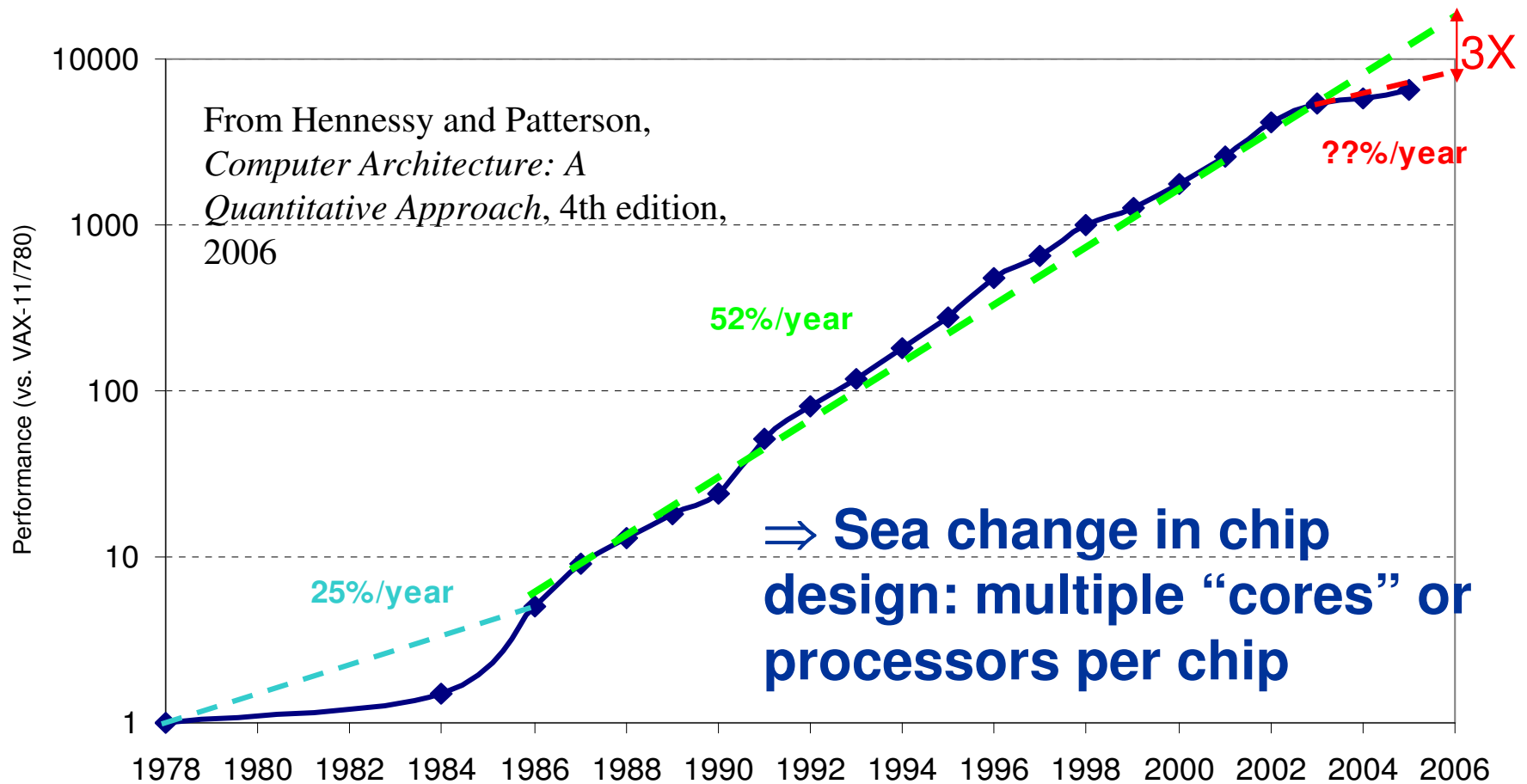
IBM Systems and Technology Group

# CMOS Microprocessor Trends, The First ~25 Years ( Good old days )

Log(Performance)

Single Thread

More active transistors, higher frequency

2005

# SPECINT



From Hennessy and Patterson,
*Computer Architecture: A
Quantitative Approach*, 4th edition,
2006

3X

??%/year

52%/year

25%/year

⇒ **Sea change in chip
design: multiple "cores" or
processors per chip**

Performance (vs. VAX-11/780)

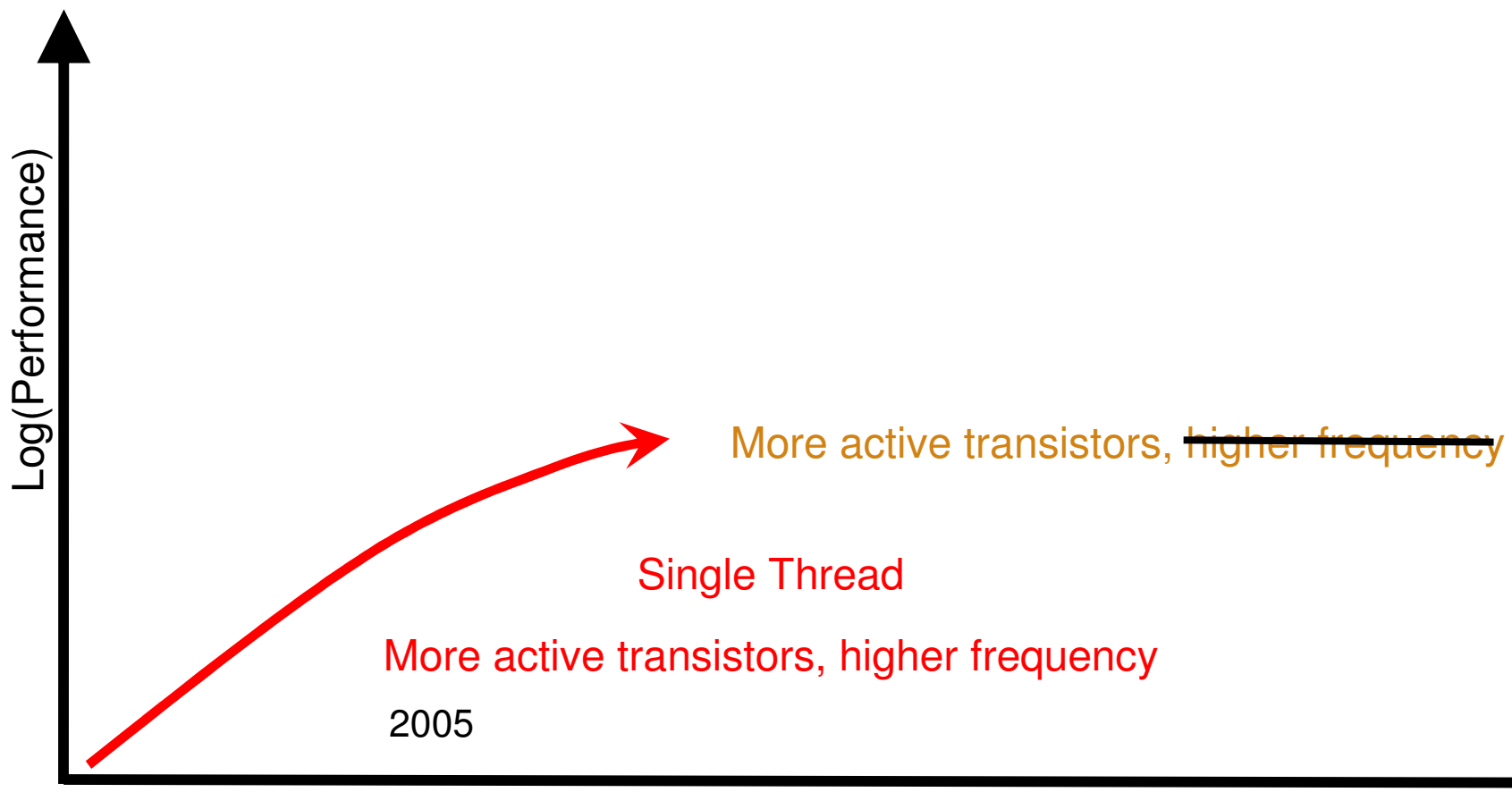**VAX          : 25%/year 1978 to 1986**
- **RISC + x86: 52%/year 1986 to 2002**
- **RISC + x86: ??%/year 2002 to present**

# Microprocessor Trends
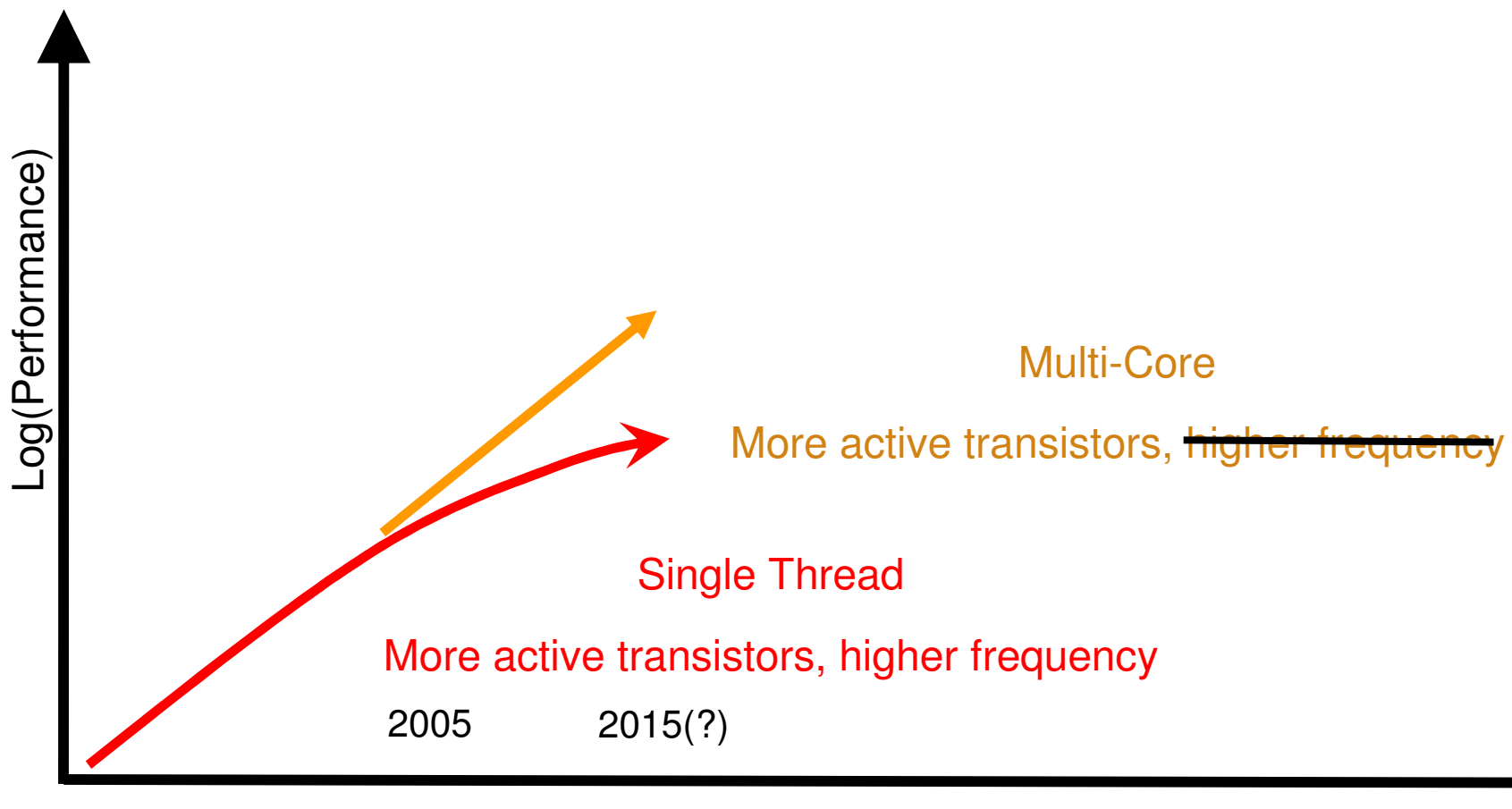
Log(Performance)

More active transistors, ~~higher frequency~~

Single Thread

More active transistors, higher frequency

2005

# CMOS Devices hit a scaling wall



Isaac e.a. IBM

# Microprocessor Trends



Log(Performance)

Multi-Core

More active transistors, ~~higher frequency~~

Single Thread

More active transistors, higher frequency

2005          2015(?)
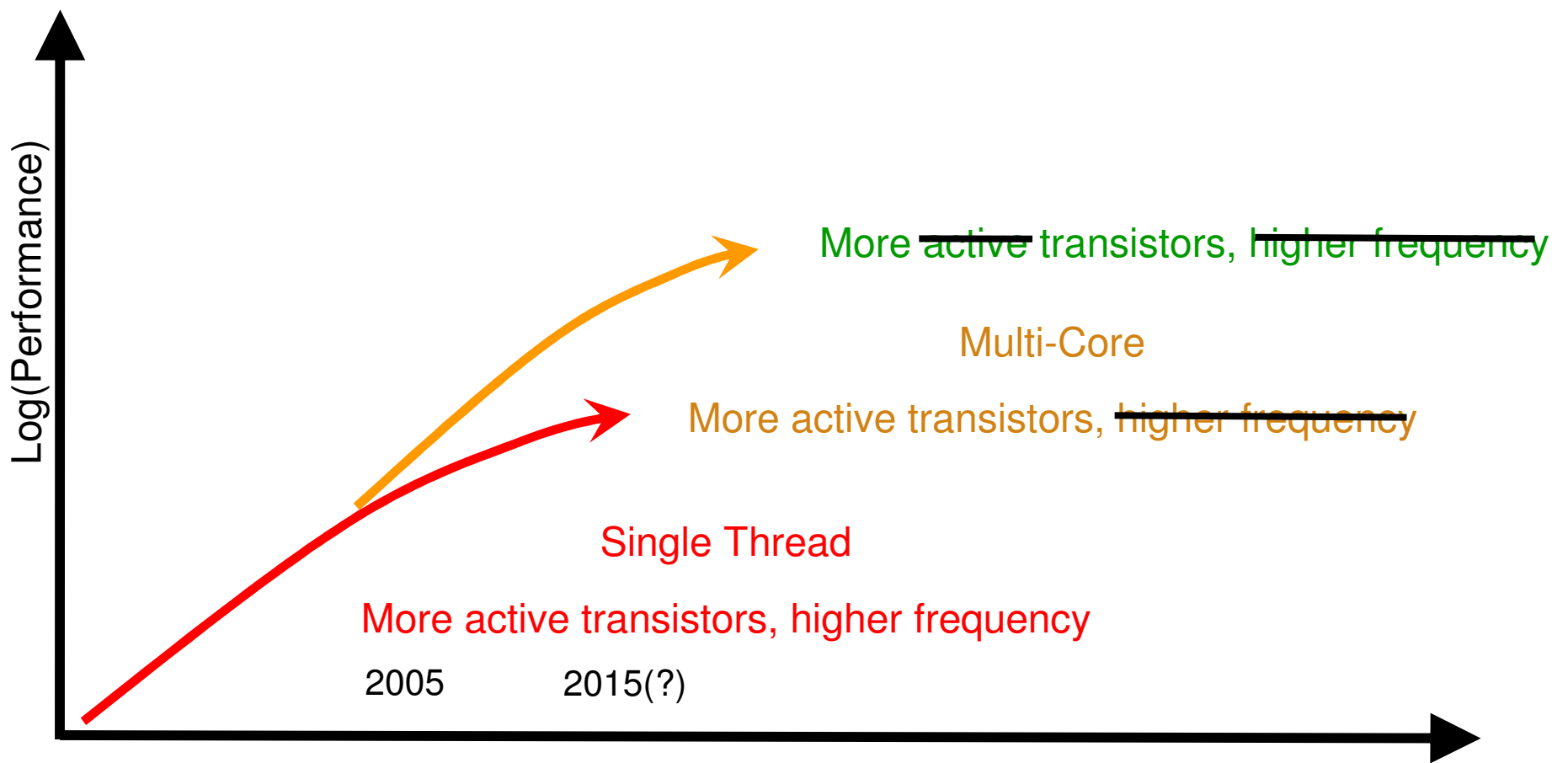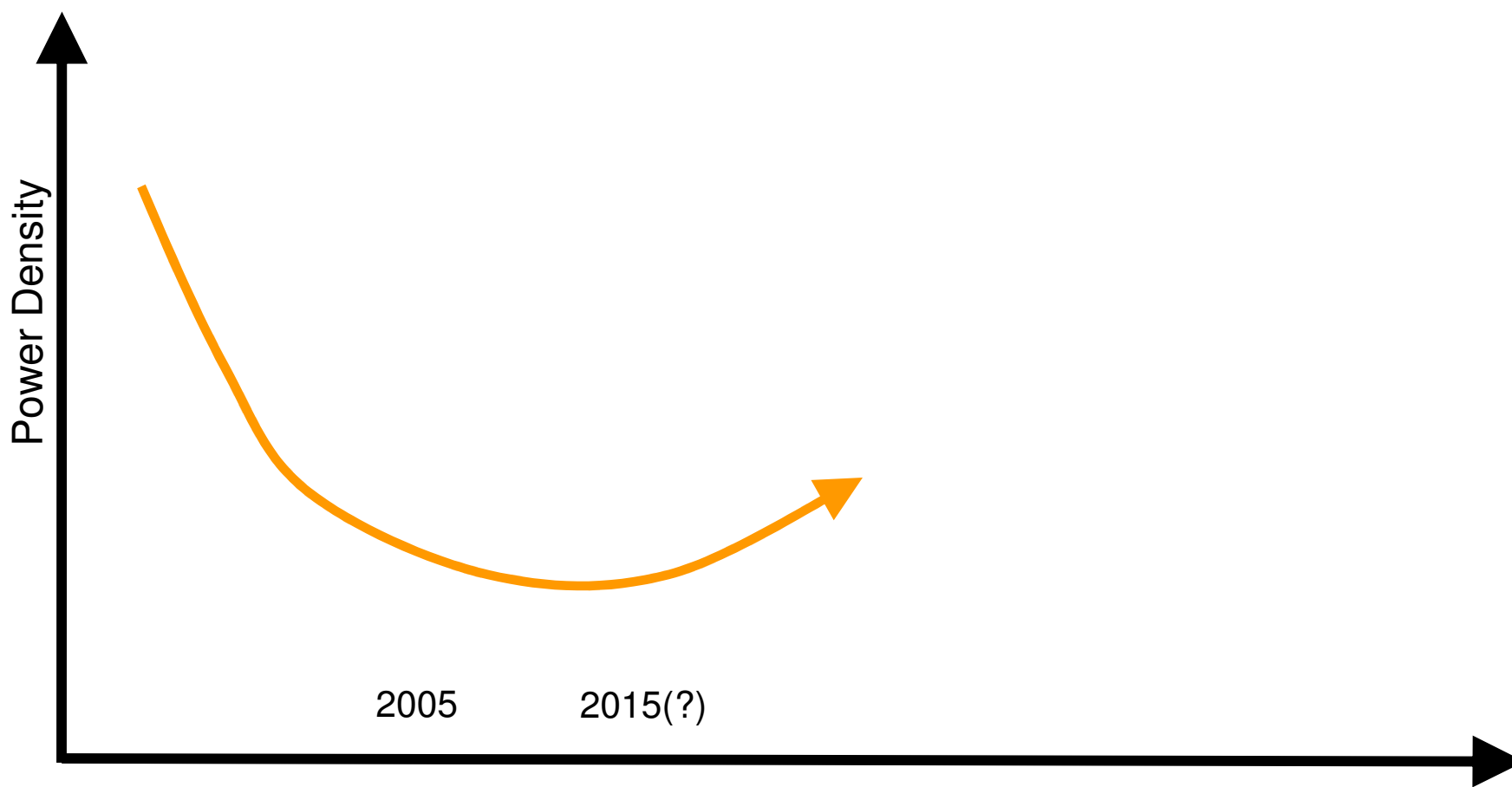
# Why are (shared memory) CMPs dominant?

- A new system delivers nearly twice the throughput performance of the previous one without application-level changes.

- Applications do not degrade in performance when ported (to a next-generation processor).
  - This is an important factor in markets where it is not possible to rewrite all applications for a new system, a common case.

- Applications benefit from more memory capacity and more memory bandwidth when ported.
  - .. even if they do not (optimally) use all the available cores.

- Even when a single application must be accelerated, large portions of code can be reused.

- Design cost is reduced, at least relative to the scenario where all available transistors are used to build a single processor.

# Microprocessor Trends



Log(Performance)

More ~~active~~ transistors, ~~higher frequency~~

Multi-Core

More active transistors, ~~higher frequency~~

Single Thread

More active transistors, higher frequency

2005          2015(?)

# Power Density at Constant Frequency



Power Density

2005    2015(?)

# Microprocessor Trends

Log(Performance)

Hybrid

More active transistors, ~~higher frequency~~

Multi-Core

More active transistors, ~~higher frequency~~

Single Thread

More active transistors, higher frequency

2005        2015(?)        2025(??)

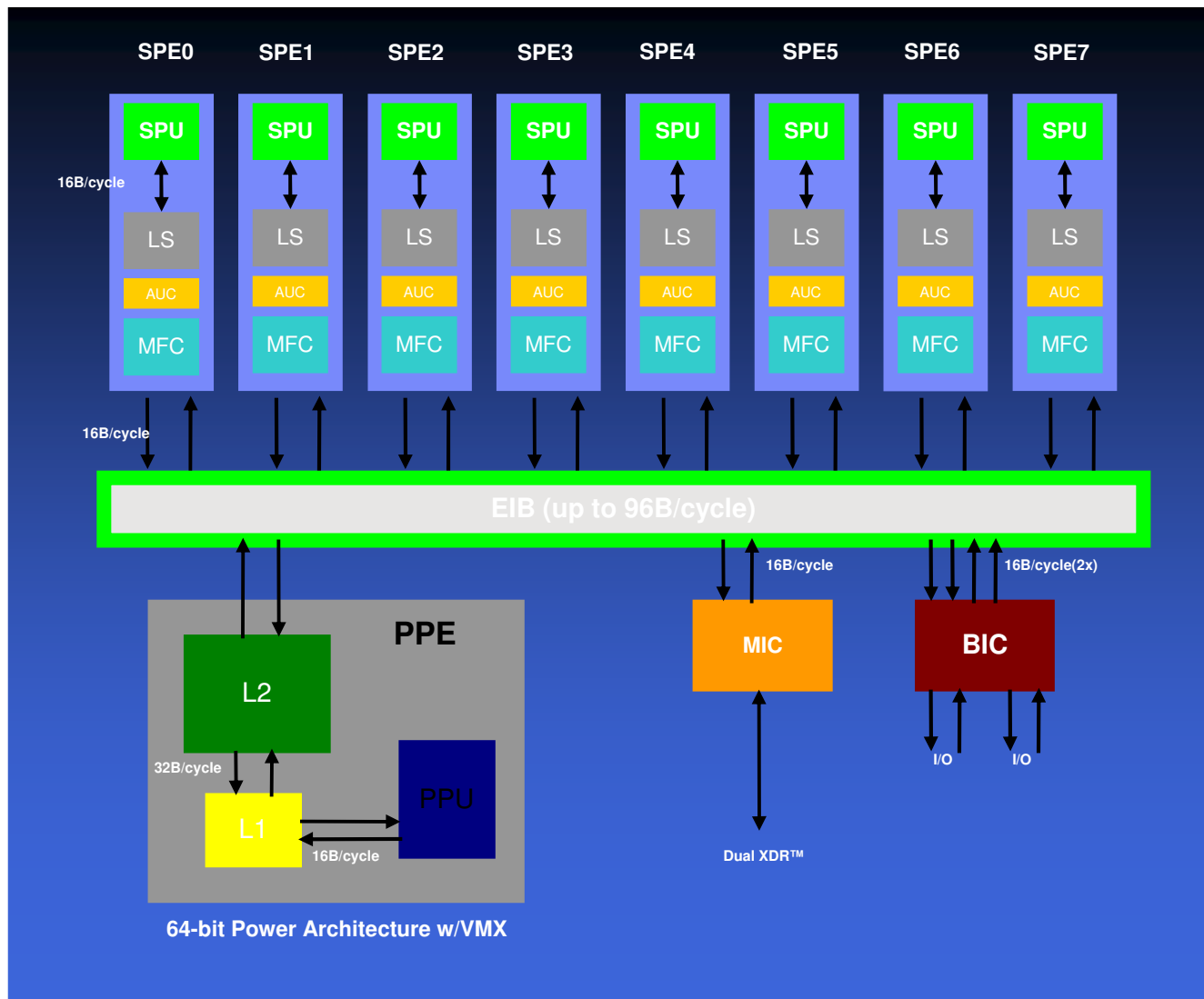# Major Sources of Efficiency in Cell Broadband Architecture

- Shopping list vs. on-demand model
- Large integrated register file
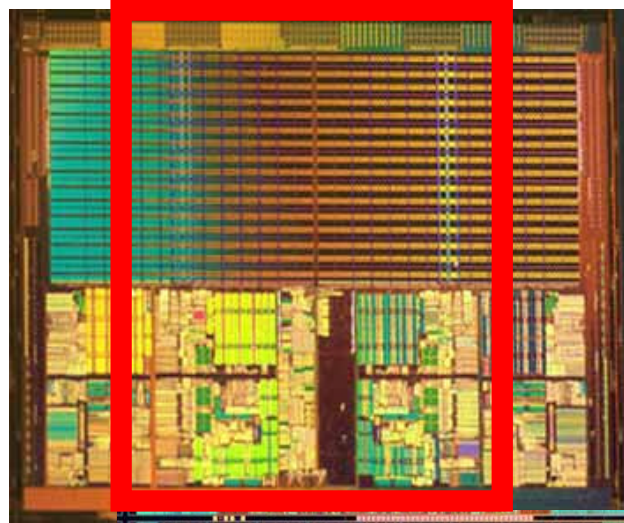- Branch hint

# Cell Broadband Engine



- Heterogeneous Multiprocessor
  - Power processor
  - Synergistic Processing Elements

- Power Processor Element (PPE)
  - general purpose
  - running full-fledged OSs
  - 2 levels of globally coherent cache

- Synergistic Proc. Element (SPE)
  - SPU optimized for computation density
  - 128 bit wide SIMD
  - Fast local memory
  - Globally coherent DMA

# Memory Managing Processor vs. Traditional General Purpose Processor



Cell BE

AMD

IBM

Intel

# IBM and its Partners are Active Users of Cell Technology

- **Three Generations of Server Blades Accompanied By 3 SDK Releases**
  - ‣ IBM QS20
  - ‣ IBM QS21
  - ‣ IBM QS22
- **Two Generations of PCIe Cell Accelerator Boards**
  - ‣ CAB ( Mercury )
  - ‣ PXCAB ( Mercury/Fixstars/Matrix Vision )
- **1U Formfactor**
  - ‣ Mercury Computer
  - ‣ TPlatforms
- **Custom Boards**
  - ‣ Hitachi Medical ( Ultrasound )
  - ‣ Other Medical and Defense
- **World's First 1 PFlop Computer**
  - ‣ LANL Roadrunner
- **Top 7 Green Systems**
  - ‣ Green 500 list

**IBM Systems**
*Simplify your IT.*

# IBM BladeCenter QS22 performance summary

The IBM BladeCenter QS22 and its IBM® PowerXCell™ 8i processor with a PPE and 8 SPEs, can perform an order of magnitude better than many traditional x86 blades when running certain applications that take advantage of the QS22's SIMD capability.

| Type | Algorithm Implementation | x86 blade / result | QS22 with IBM PowerXCell 8i 3.2 GHz processor(s) / result | Comparison Factor |
|---|---|---|---|---|
| High Performance Computing (HPC) | Matrix Multiplication (S.P.) | 86 blade (2.66GHz Quad-Core Intel® X5355) / 77 GFlops | 8 SPEs / 203 GFlops | x<br>2.6x |
| | LINPACK (S.P.) | x86 blade (2.66GHz Quad-Core Intel X5355) / 73 GFlops | 8 SPEs / 164 GFlops | 2.2x |
| | Matrix Multiplication (D.P.) | x86 blade (2.66GHz Quad-Core Intel X5355) / 38 GFlops | 8 SPEs / 101 GFlops | 2.6x |
| | LINPACK (D.P.) | x86 blade (2.66GHz Quad-Core Intel X5355) / 36 GFlops | 8 SPEs / 84.8 GFlops | .3x |
| | 3-step 2D PFAFFT | x86 blade (3.0 GHz Dual-Core Intel X5160 x2) / 16 - 687 seconds | 16 SPEs / 6.6 - 89 seconds | 2.4-7.7x |
| Medical / HCLS | HMMer | x86 blade (3.0 GHz Dual-Core Intel X5160) / 428 sec | 8 SPEs / 34.4 sec | 12.4x |
| Financial Services Sector (FSS) | merican Option using Binomial Tree | 86 blade (2.33 GHz Quad-Core Intel E5345) / 19K Options per second | 8 SPEs / 107K Options per second | 5.6x |
| | ollateralized Debt Obligation (CDO) | x86 blade (2.8 GHz Quad-Core Intel E5440) / 28 TSps | 8 SPEs / 211 TSps | 7.5x |

**The source for all data is IBM internal benchmark testing as of April 15, 2008.  Different applications implementing these algorithms may affect performance results.  These results were derived using particular hardware and software configurations; differences in hardware and software configurations may affect performance results.**

Notes: refer to "Notes on Benchmarks and Values" chart;  S.P.: Single Precision; D.P.: Double Precision; GFlops:Giga Floating point operations per second; seconds: Elapsed time in seconds; second;Gbps=Gigabits per second; PPE: Power Processing Element; SPE: Synergistic Processing Element

**IBM Systems**
*Simplify your IT.*

QS22 with IBM PowerXCell 8i

| Type | Algorithm implementation | x86 blade / result | 3.2 GHz processor(s) / result | Comparison Factor |
|---|---|---|---|---|
| Financial Services Sector (FSS) cont. | uropean Options using Black-Scholes (D.P.) | x86 blade (2.66 GHz Quad-Core Intel X5355) D.P. 35 MBOPS | 8 SPEs / 125 MBOPS | 3.5x |
| | European Options using Monte-Carlo | x86 blade (2.33 GHz Quad-Core Intel E5345) S.P. 210 MSps D.P. 65-122 MSps | 8 SPEs / S.P. 1300 MSps D.P. 291-325 MSps | S.P. 6.1x D.P. 2.6-4.4x |
| Linear Algebra Libraries | BLAS routines | x86 blade (2.33 GHz Quad-Core Intel E5345 x 2) DDOT: 0.37 GFlops DAXPY: 0.27 GFlops DTRMM: 41 GFlops | 6 SPEs / DOT: 1.9 GFlops AXPY: 1.4 GFlops TRMM: 123 GFlops | DDOT: 5.1x DAXPY: 5.1X DTRMM: 3.0X |
| | LAPACK routines | x86 blade (2.33 GHz Quad-Core Intel E5345 x 2) DGETRF: 28 GFlops DPOTRF: 31.7 GFlops | 16 SPEs / DGETRF: 105 GFlops DPOTRF: 140 GFlops | DGETRF: 3.7X DPOTRF: 4.4X |

**QS22 results where no comparison data was gathered:**

| Type | Algorithm Implementation | QS22 with IBM PowerXCell 8i 3.2 GHz processor(s) result |
|---|---|---|
| High Performance Computing (HPC) | SCAMPI Network Intrusion Detection | 16 SPEs / 13-16 Gbps |
| Digital Media | IBM iRT Demo of Boeing 777 | 112 SPEs (14 QS22's in single IBM BladeCenter) / More than 5 frames per second for 25 GB size model containing 300M triangles |
| Medical / HCLS | Rigid Tissue Image Registration | 16 SPEs / 158 seconds for 94 images |

Notes: Refer to "Notes on Benchmarks and Values" chart; MBOPS: Million Blackscholes operations per sec; S.P.: Single Precision; D.P. Double Precision; MSps: Million Simulations per second; GFlops: Giga Floating Operations per second; SPE: Synergistic Processing Element

# Optimization of Sparse Matrix-Vector Multiplication on Emerging Multicore Platforms

Samuel Williams[*][†], Leonid Oliker[*], Richard Vuduc[§], John Shalf[*], Katherine Yelick[*][†], James Demmel[†]

[*]CRD/NERSC, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA
[†]Computer Science Division, University of California at Berkeley, Berkeley, CA 94720, USA
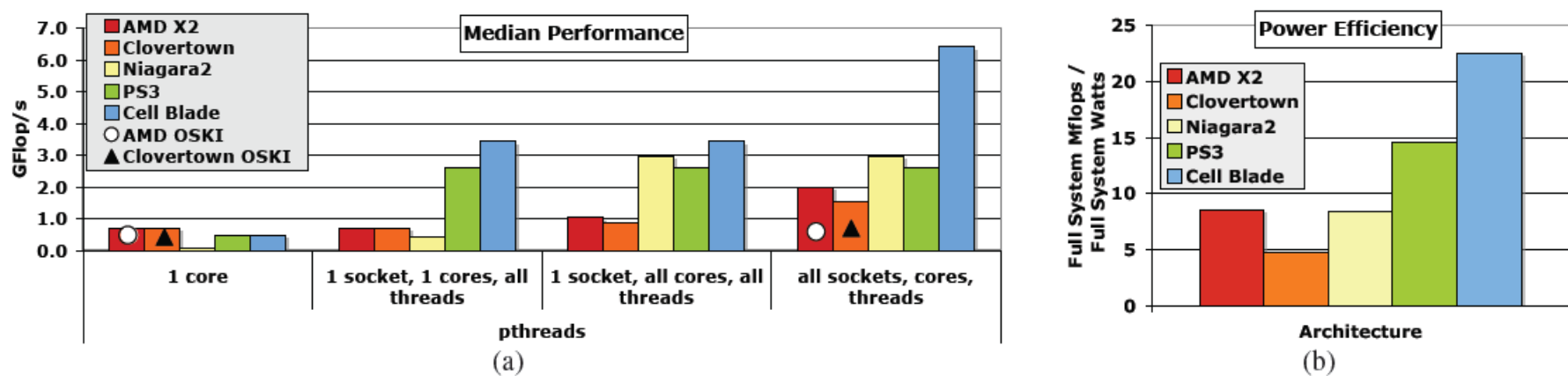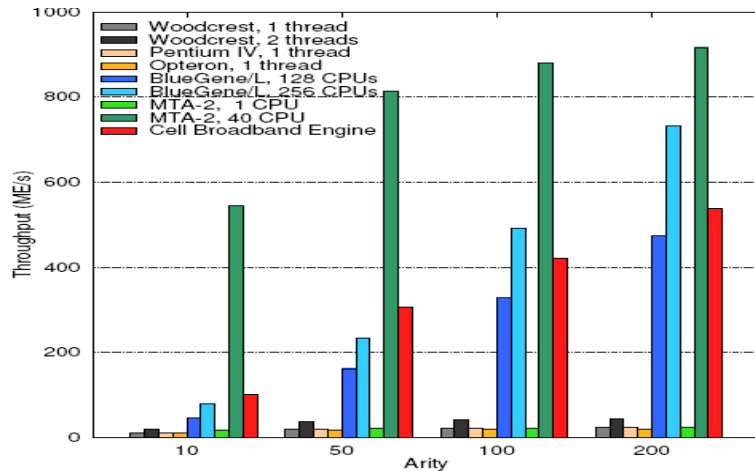[§]CASC, Lawrence Livermore National Laboratory, Livermore, CA 94551, USA

Figure 5: Architectural comparison of the median matrix performance showing (a) GFlop/s rates of OSKI and optimized SpMV on single-core, full socket, and full system and (b) relative power efficiency computed as total full system Mflop/s divided by sustained full system Watts (see Table 1).
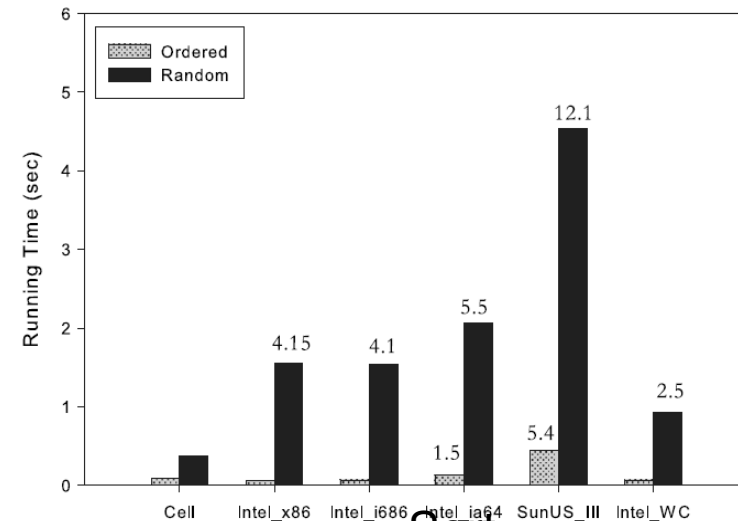
# Current Cell: Integer Workloads

Breadth-First Search
Villa, Scarpazza, Petrini, Peinador
IPDPS 2007



**a** Comparison of List ranking on Cell with other Single Processors for list of size 8 million nodes



## Mapreduce
### Sangkaralingam, De Kruijf, Oct. 2007

Sort
Gedik, Bordawekar, Yu (IBM)

| Application Name | Application Type | Lines of Code | | Speedup vs. Core2 | | | BIPS | | |
|---|---|---|---|---|---|---|---|---|---|
| | | MapReduce | Serial | 1-SPE | 8-SPEs | 8-SPE Ideal | 1-SPE | 8-SPEs | 8-SPE Ideal |
| histogram | partition-dominated | 345 | 216 | 0.16 | 0.15 | 2.44 | 1.56 | 1.51 | 24.49 |
| kmeans | partition-dominated | 324 | 318 | 0.91 | 3.00 | 6.92 | 2.08 | 7.35 | 17.01 |
| linearRegression | map-dominated | 279 | 114 | 0.34 | 2.59 | 2.67 | 1.47 | 11.32 | 11.70 |
| wordCount | partition-dominated | 226 | 324 | 0.87 | 0.96 | 10.26 | 1.52 | 1.74 | 18.64 |
| NAS_EP | map-dominated | 264 | 112 | 1.08 | 8.62 | 8.62 | 2.00 | 15.93 | 15.95 |
| distributedSort | sort-dominated | 171 | 93[c] | 0.41 | 0.76 | 5.48 | 1.28 | 2.38 | 17.15 |

Table 3: Out-of-core sort performance (in secs)

| # items | 16 SPEs bitonic | 3.2GHz Xeon quick | 3.2GHz Xeon quick 2-core | PPE quick |
|---|---|---|---|---|
| 1M | 0.0098 | 0.1813 | 0.098589 | 0.4333 |
| 2M | 0.0234 | 0.3794 | 0.205728 | 0.9072 |
| 4M | 0.0569 | 0.7941 | 0.429499 | 1.9574 |
| 8M | 0.1372 | 1.6704 | 0.895168 | 4.0746 |
| 16M | 0.3172 | 3.4673 | 1.863354 | 8.4577 |
| 32M | 0.7461 | 7.1751 | 3.863495 | 18.3882 |
| 64M | 1.7703 | 14.8731 | 7.946356 | 38.7473 |
| 128M | 4.0991 | 30.0481 | 16.165578 | 79.9971 |

IBM Systems
*Simplify your IT.*

# Microprocessor Trends



Log(Performance)

More active transistors, higher frequency

Hybrid

More active transistors, higher frequency

Multi-Core

More active transistors, higher frequency

Single Thread

More active transistors, higher frequency

2005        2015(?)        2025(??)        2035(???)

# Microprocessor Trends

Log(Performance)

**Special Purpose ( ASIC )**
More active transistors, higher frequency

**Hybrid**
More active transistors, higher frequency

**Multi-Core**
More active transistors, higher frequency

**Single Thread**

More active transistors, higher frequency

2005          2015(?)          2025(??)          2035(???)

# Five Decades of Innovations

Airgap

**S/360 Model 67 first virtualized machine**

Thermal conduction cooling technology

CMOS processors

Modular refrigeration cooling technology

High-k metal gates

POWER6

IBM Energy Efficiency Institute, Austin, TX

**1960s -1970s**

**mid-1990s**

**2000s**

**1980s**

VM virtualization

late-1990s

Cell BE processor

eDRAM

Air / liquid hybrid cooling technology

Flat plate conduction cooling technology

Copper chip

3D chip stacking

**IBM Systems** *Simplify your IT.*

# Performance and Productivity Challenges require a Multi-Dimensional Approach

| Highly Productive Systems | Highly Scalable Multi-core Systems | Hybrid Systems |
|---|---|---|
| **POWER** | | |

**Comprehensive (Holistic) System Innovation & Optimization**

# HPC Cluster Directions



**ExaScale**

**ExaF**

Accelerators

**Targeted Configurability**

BG/Q    Accelerators

**Performance**

**PF**

Roadrunner

PERCS Systems

Accelerators

**Capability Machines**

BG/P

**Extended Configurability**

Accelerators

Power, AIX/ Linux

Accelerators

**Capacity Clusters**

Accelerators

Linux Clusters  Power, x86-64,
Less Demanding Communication

**2007    2008    2009/2010    2011    2012    2013    2018-19**

**IBM Systems**
*Simplify your IT.*

# Next Era of Innovation – Hybrid Computing
# The Next Bold Step in Innovation & Integration

**Symmetric Multiprocessing Era**

**Hybrid Computing Era**
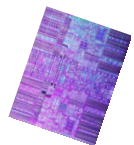
*Today*

*pNext 1.0*

*pNext 2.0*

**p6**

**p7**

Traditional

**Cell**

Throughput

Computational

**BlueGene**

Technology Out

**Driven by cores/threads**

Market In

**Driven by workload consolidation**

**IBM Systems** *Simplify your IT.*

# Programming Models: Architecture

## User View

| PGAS CAF /X10/ UPC | Shmem GSM | MPI Fortran C |
|---|---|---|

| Logical Single Address Space | Logical Multiple Address Spaces |
|---|---|

**Logical View**

| Multiple Machine Address Spaces |
|---|

**Cluster Level**

| Homogeneous Cores | Heterogeneous |
|---|---|
| BG | Power/PERC | Roadrunner |
| Open MP, SM-MPI | Open MP, SM-MPI | Open MP, OpenCL |
| HW Cache | HW Cache | Software Cache |

**Node Level**

# Cell/B.E. Soft I-Cache Summary

- Up to ½ GB of code
- Normal tool-chain flow
  - ▸ No detailed knowledge required on the part of the developer.
- Use self-modifying code (mini-JIT) to have branches go directly to their targets when they are in cache – no overhead in hit case.
  - ▸ **Less than 10% total runtime penalty for running in small caches. Still working to improve.**
  - ▸ Verified on QS22 & PXCAB hardware.
- Support code out-side of cache structure
- 'Small' changes to ABI – good operability with old source.
  - ▸ New virtual address space for code
    - – 32 bit function pointers
    - – Indirects require tag check

**.C**

**compiler**

Divide code with branch always into Blocks smaller than the line size & annotate branches with importance+ Return stack information

**.S**

**assembler**

**.O**

Linker analyzes whole program to Determine cache layout that minimizes cache conflicts on Important paths
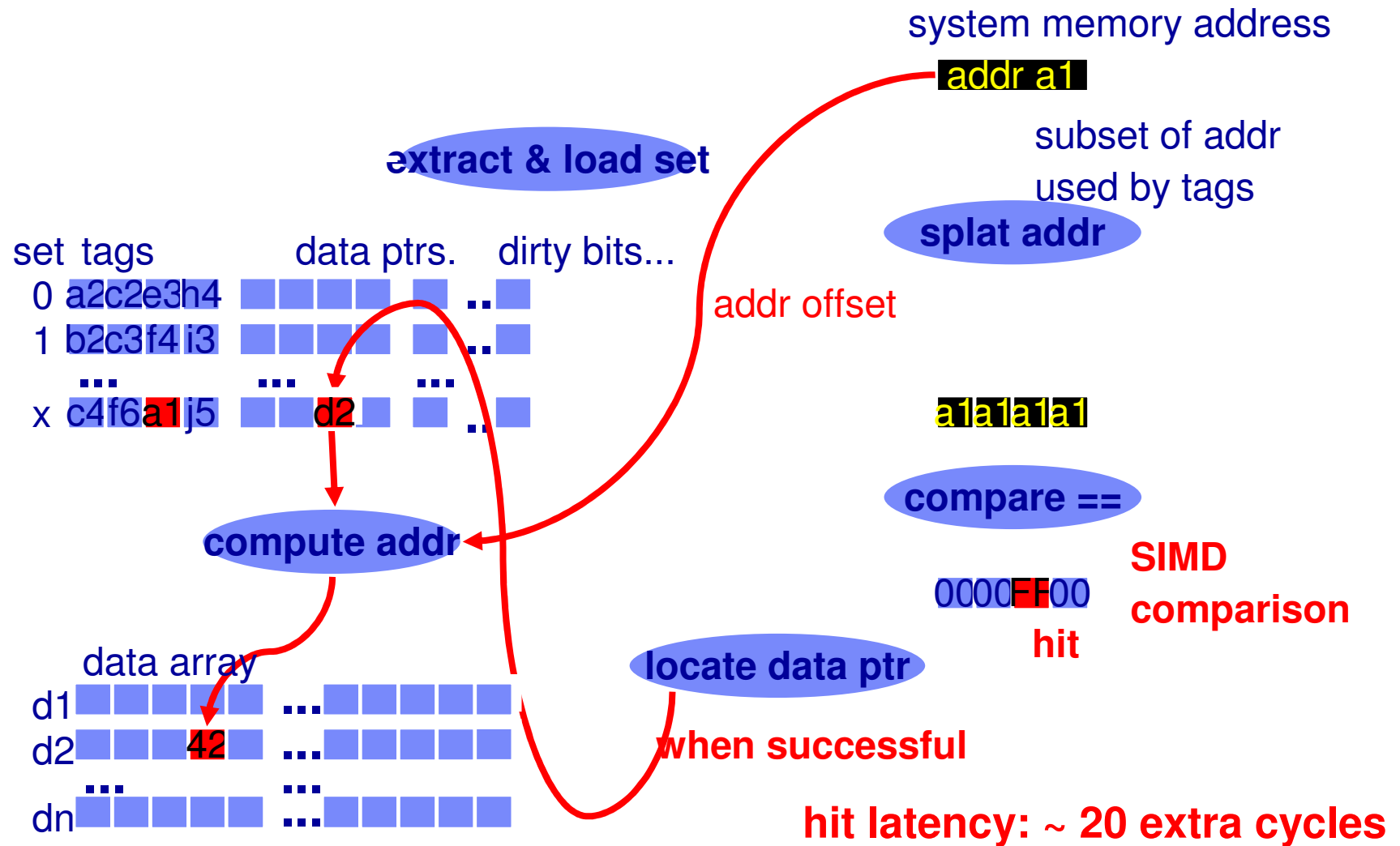
**linker**

**Runtime System**

**Exe/lib**

B. Flachs e.a.. IBM

# Software Data Cache ( XL Compiler )

- Works just like a hardware cache, but implemented in software

  ‣ Loads/stores replaced with software cache lookup instructions

  ‣ Miss handler invoked for a cache miss

    • Brings in the missing cache line, evicts an existing cache line if necessary

  ‣ 128B cache line, 4-way associative. Cache size configurable with command line option

- Coherence among threads

  ‣ One cache line may be shared by multiple SPE threads – cannot naïvely evict whole cache line

  ‣ Dirty bits to record modified data (in unit of bytes)

  ‣ Atomic updates based on dirty bits to evict a cache line

- Pros/Cons

  ‣ Uniform solution for all kinds of references

  ‣ Exploit data reuse dynamically

  ‣ **Overhead ( Unlike SW-ICache, SW-DCache generally not competitive with hardware )**

M. Mendel, K. O'Brien, e.a., IBM

# Software Data Cache Access



system memory address

addr a1

subset of addr
used by tags

**extract & load set**

**splat addr**

set  tags                data ptrs.   dirty bits...

0  a2c2e3h4  ▪▪▪▪ ▪ ▪..▪
1  b2c3f4i3  ▪▪▪▪ ▪ ▪..▪
   ...        ...        ...
x  c4f6a1j5  ▪▪▪ d2 ▪ ▪..▪

addr offset

a1a1a1a1

**compare ==**

**compute addr**

0000FF00

**SIMD comparison**

hit

data array

**locate data ptr**

d1  ▪▪▪▪▪ ▪▪... ▪▪▪▪▪
d2  ▪▪▪ 42 ▪ ▪▪... ▪▪▪▪▪
   ...        ...
dn  ▪▪▪▪▪ ▪▪... ▪▪▪▪▪

**when successful**

**hit latency: ~ 20 extra cycles**

# DMA Tiling ( XL Compiler )

- **Handles regular data accesses to shared memory by compiler**
  - ▸ Buffers in SPE local memory are controlled by compiler
  - ▸ Calls to allocate and free buffers are inserted
  - ▸ DMA operations are inserted
  - ▸ References to global variables are replaced by direct references to the local buffer

- **Pros/Cons**
  - ▸ Much less overhead: no lookup, more control on DMA
  - ▸ Compile time decision to use DMA tiling
    - • Sometimes not possible
    - • Sometimes not optimal
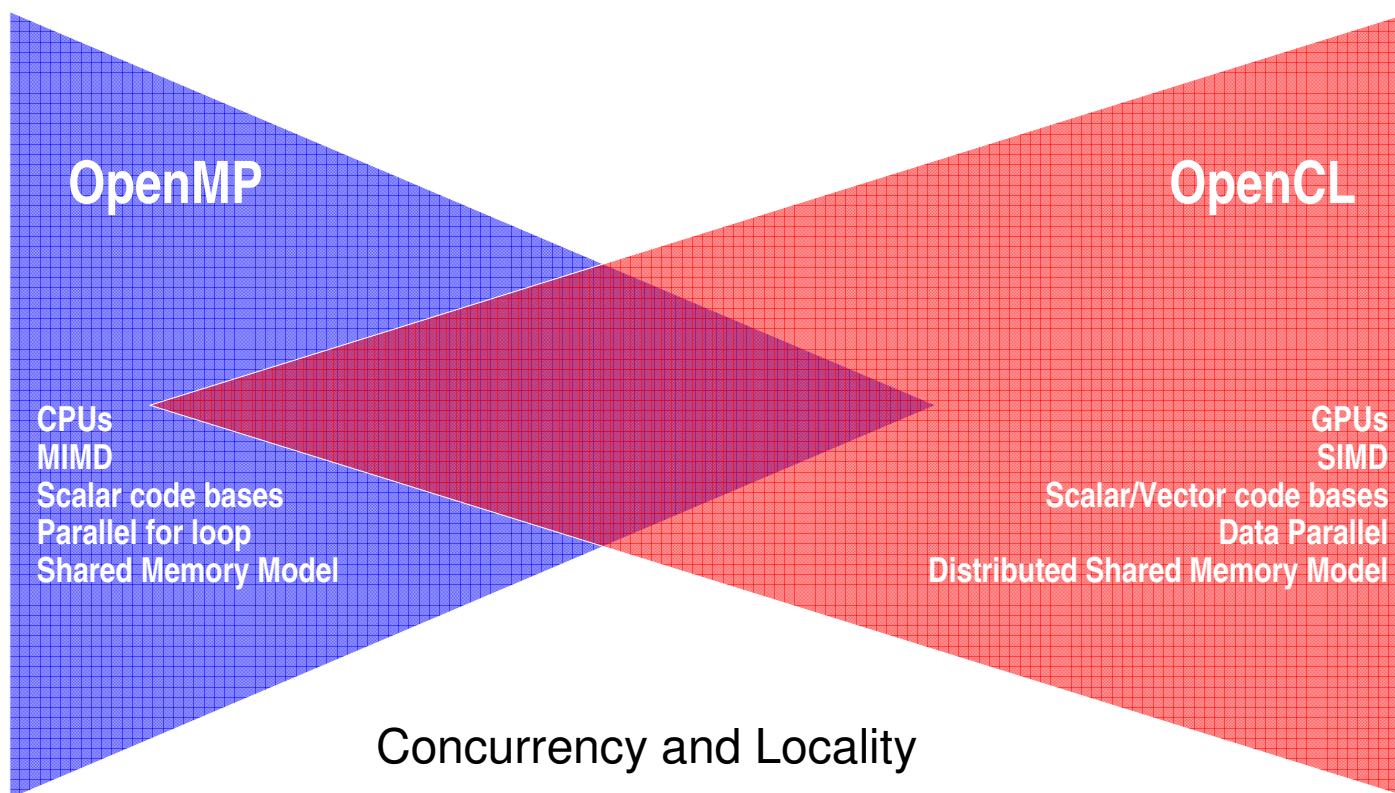
```
for (i=0; i<N; i++) {
    A[i] = B[i]*C[i]
}
```

→

```
for (ii=0; i<N; i+=bf) {
        read part B into B';
        read part C into C';
        for (i=ii; i < min(ii+bf, N); i++) {
                A'[i]=B'[i]*C'[i];
        }
        write A' back to A;
}
```

# Two Standards for Programming the Node

- Two standards evolving from different sides of the market

**OpenMP**

**OpenCL**

CPUs
MIMD
Scalar code bases
Parallel for loop
Shared Memory Model

GPUs
SIMD
Scalar/Vector code bases
Data Parallel
Distributed Shared Memory Model

Concurrency and Locality

**IBM Systems**
*Simplify your IT.*

# Cell Broadband Engine

- Unified host and device memory
  - Zero copies between them



**SPE 1**

WorkItem 1

Compute Unit 1

Local Memory

Private Memory

Local Store

**SPE N**

WorkItem N

Compute Unit N

Local Memory

Private Memory

Local Store

**PPE**

Host Device

Compute Device

Host Memory

Device Global Memory

System Memory

# Cell/B.E. observations

- TASKS!
  - By programmer
  - In runtime
  - In language
  - In acceleration paradigm

  Nice because:
  - Scalable
  - No load-balancing concerns
  - Much less opportunity for difficult MP-issues

# Summary

- ## Technology limits drive fundamental change:
  - ▸ First multi-core, then hybrid and eventually special-purpose?
  - ▸ Cell an early example of hybrid

- ## What is next:
  - ▸ Continued Focus on Efficiency
  - ▸ Increasing Focus on Standards-Based Programming
    - – Software ICache & Software DCache for Cell/B.E.
    - – OpenMP & OpenCL for Cell/B.E. and other processors
    - – …
  - ▸ Increasing Focus on Ease of Use
    - – Make accelerators "invisible" for most customers
    - – Commercial applications, not just HPC
    - – Not an easy thing to do

  - ▸ Continue to Broaden Application Reach for Cell and Hybrid Systems

**IBM Systems**
*Simplify your IT.*

# Special notices

This document was developed for IBM offerings in the United States as of the date of publication.  IBM may not make these offerings available in other countries, and the information is subject to change without notice. Consult your local IBM business contact for information on the IBM offerings available in your area.

Information in this document concerning non-IBM products was obtained from the suppliers of these products or other public sources.  Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

IBM may have patents or pending patent applications covering subject matter in this document.  The furnishing of this document does not give you any license to these patents.  Send license inquires, in writing, to IBM Director of Licensing, IBM Corporation, New Castle Drive, Armonk, NY 10504-1785 USA.

All statements regarding IBM future direction and intent are subject to change or withdrawal without notice, and represent goals and objectives only.

The information contained in this document has not been submitted to any formal IBM test and is provided "AS IS" with no warranties or guarantees either expressed or implied.

All examples cited or described in this document are presented as illustrations of  the manner in which some IBM products can be used and the results that may be achieved.  Actual environmental costs and performance characteristics will vary depending on individual client configurations and conditions.

IBM Global Financing offerings are provided through IBM Credit Corporation in the United States and other IBM subsidiaries and divisions worldwide to qualified commercial and government clients.  Rates are based on a client's credit rating, financing terms, offering type, equipment type and options, and may vary by country.  Other restrictions may apply.  Rates and offerings are subject to change, extension or withdrawal without notice.

IBM is not responsible for printing errors in this document that result in pricing or information inaccuracies.

All prices shown are IBM's United States suggested list prices and are subject to change without notice; reseller prices may vary.

IBM hardware products are manufactured from new parts, or new and serviceable used parts. Regardless, our warranty terms apply.

Any performance data contained in this document was determined in a controlled environment.  Actual results may vary significantly and are dependent on many factors including system hardware configuration and software design and configuration.  Some measurements quoted in this document may have been made on development-level systems.  There is no guarantee these measurements will be the same on generally-available systems.  Some measurements quoted in this document may have been estimated through extrapolation.  Users of this document should verify the applicable data for their specific environment.

# Special notices (cont.)

The following terms are registered trademarks of International Business Machines Corporation in the United States and/or other countries: AIX, AIX/L, AIX/L(logo), alphaWorks, AS/400, BladeCenter, Blue Gene, Blue Lightning, C Set++, CICS, CICS/6000, ClusterProven, CT/2, DataHub, DataJoiner, DB2, DEEP BLUE, developerWorks, DirectTalk, Domino, DYNIX, DYNIX/ptx, e business(logo), e(logo)business, e(logo)server, Enterprise Storage Server, ESCON, FlashCopy, GDDM, i5/OS, IBM, IBM(logo), ibm.com, IBM Business Partner (logo), Informix, IntelliStation, IQ-Link, LANStreamer, LoadLeveler, Lotus, Lotus Notes, Lotusphere, Magstar, MediaStreamer, Micro Channel, MQSeries, Net.Data, Netfinity, NetView, Network Station, Notes, NUMA-Q, OpenPower, Operating System/2, Operating System/400, OS/2, OS/390, OS/400, Parallel Sysplex, PartnerLink, PartnerWorld, Passport Advantage, POWERparallel, Power PC 603, Power PC 604, PowerPC, PowerPC(logo), Predictive Failure Analysis,  pSeries, PTX, ptx/ADMIN, RETAIN, RISC System/6000, RS/6000, RT Personal Computer, S/390, Scalable POWERparallel Systems, SecureWay, Sequent, ServerProven, SpaceBall, System/390, The Engines of e-business, THINK, Tivoli, Tivoli(logo), Tivoli Management Environment, Tivoli Ready(logo), TME, TotalStorage, TURBOWAYS, VisualAge, WebSphere, xSeries, z/OS, zSeries.

The following terms are trademarks of International Business Machines Corporation in the United States and/or other countries: Advanced Micro-Partitioning, AIX 5L, AIX PVMe, AS/400e, Chiphopper, Chipkill, Cloudscape, DB2 OLAP Server, DB2 Universal Database, DFDSM, DFSORT, DS4000, DS6000, DS8000, e-business(logo), e-business on demand, eServer, Express Middleware, Express Portfolio, Express Servers, Express Servers and Storage, General Purpose File System, GigaProcessor, GPFS, HACMP, HACMP/6000, IBM TotalStorage Proven, IBMLink, IMS, Intelligent Miner, iSeries, Micro-Partitioning, NUMACenter, On Demand Business logo, POWER, PowerExecutive, Power Architecture, Power Everywhere, Power Family, Power PC, PowerPC Architecture, PowerPC 603, PowerPC 603e, PowerPC 604, PowerPC 750, POWER2, POWER2 Architecture, POWER3, POWER4, POWER4+, POWER5, POWER5+, POWER6, POWER6+, pure XML, Redbooks, Sequent (logo), SequentLINK, Server Advantage, ServeRAID, Service Director, SmoothStart, SP, System i, System i5, System p, System p5, System Storage, System z, System z9, S/390 Parallel Enterprise Server, Tivoli Enterprise, TME 10, TotalStorage Proven, Ultramedia, VideoCharger, Virtualization Engine, Visualization Data Explorer, X-Architecture, z/Architecture, z/9.

A full list of U.S. trademarks owned by IBM may be found at: http://www.ibm.com/legal/copytrade.shtml.
The Power Architecture and Power.org wordmarks and the Power and Power.org logos and related marks are trademarks and service marks licensed by Power.org.
UNIX is a registered trademark of The Open Group in the United States, other countries or both.
Linux is a trademark of Linus Torvalds in the United States, other countries or both.
Microsoft, Windows, Windows NT and the Windows logo are registered trademarks of Microsoft Corporation in the United States, other countries or both.
Intel, Itanium, Pentium are registered tradas and Xeon is a trademark of Intel Corporation or its subsidiaries in the United States, other countries or both.
AMD Opteron is a trademark of Advanced Micro Devices, Inc.
Java and all Java-based trademarks and logos are trademarks of Sun Microsystems, Inc. in the United States, other countries or both.
TPC-C and TPC-H are trademarks of the Transaction Performance Processing Council (TPPC).
SPECint, SPECfp, SPECjbb, SPECweb, SPECjAppServer, SPEC OMP, SPECviewperf, SPECapc, SPEChpc, SPECjvm, SPECmail, SPECimap and SPECsfs are trademarks of the Standard Performance Evaluation Corp (SPEC).
NetBench is a registered trademark of Ziff Davis Media in the United States, other countries or both.
AltiVec is a trademark of Freescale Semiconductor, Inc.
Cell Broadband Engine is a trademark of Sony Computer Entertainment Inc.
Other company, product and service names may be trademarks or service marks of others.

# Panel

- Will a typical CS graduate be able to program mainstream, projected many-core architectures?
  - ▸ Yes, but not efficiently.
- Is there a road to portability between different types of many-core architectures?
  - ▸ OpenMP & OpenCL ( Many-core on-chip )
- If not, should the major vendors look for other, perhaps more innovative, approaches to (highly) parallel many-core architectures?
  - ▸ More innovative = less portable?
- What characteristics should such many-core architectures have?
  - ▸ (Chip) Hardware model should be based on shared memory but able to leverage locality and predictability (reuse/prefetch-ability) for added performance.
- Can programming models, parallel languages, libraries, and other software help?
  - ▸ Enhance task model in OpenMP and OpenCL. Better runtimes!
- Is parallel processing research on track?
  - ▸ Not much fundamental treatment of locality.
- What will the typical CS student need in the coming years?
  - ▸ A much more fundamental understanding of how algorithms map to hardware.