


## **Towards Highly Scalable Clusters for Crash**

Dr. Achim Bömelburg,


IBM Germany



IBM Deep Computing Team


## Towards Highly Scalable Clusters for Crash

Dr. Achim Bömelburg, IBM Germany



Spetember 2007

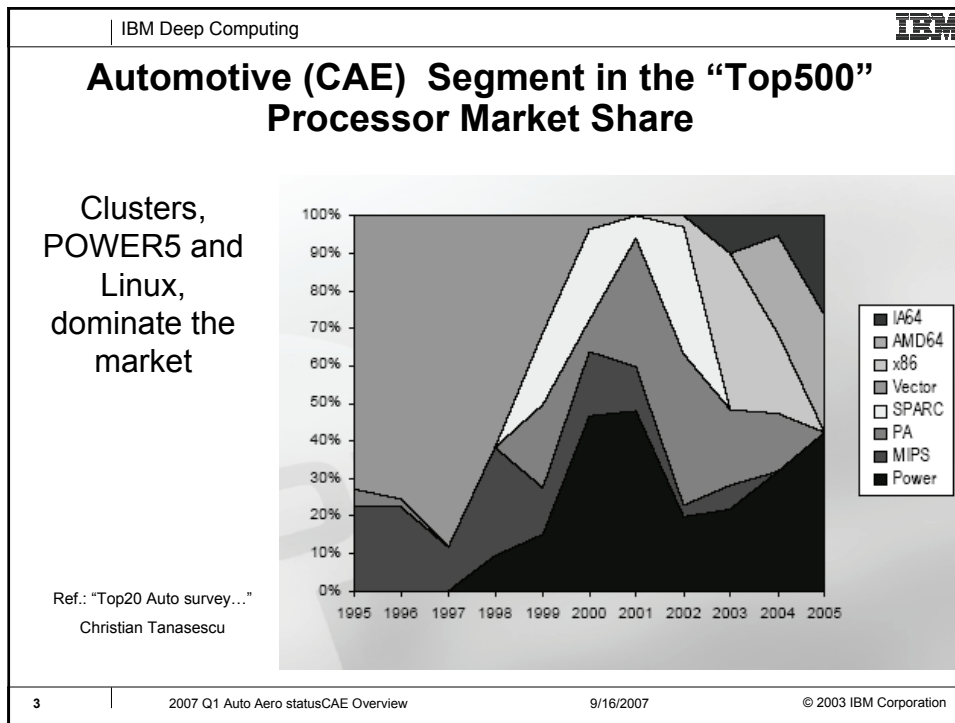
2007 Q1 Auto Aero status9/16/2007© 2006 IBM Corporation


IBM Deep Computing Team

## Contents

- **Industry Trends**
  - Auto/HPC trends
- **HPC Product Update**
  - XEON 5100 series
  - Opteron “rev F”
- **IBM “value add”**
  - GPFS
  - vMIO
  - 10M element crash
- **General Discussion**

22007 Q1 Auto Aero status9/16/2007© 2006 IBM Corporation

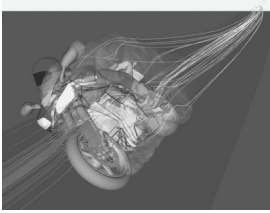


IBM Deep Computing 

## Leadership is a result of good decisions and solid execution

- **Investment in POWER processor line (1999)**
  - Increased commitment to HPC market
  - Leverage technology across HPC and commercial
- **Introduction of POWER4 (2001)**
  - First "dual-core" processor
  - 2x performance advantage over competition
- **Early investment/commitment in Linux (2001)**
  - Expertise in place for Auto industry transition.
- **First "tier 1" vendor with Opteron (2004)**
  - Workstation, servers and blades
  - Did not invest in Itanium
- **Investment in "blade" technology**
  - Offerings for POWER, Xeon, and Opteron

Image courtesy of CEI



**IBM is a driver ...  
others are  
passengers**

4 | 2007 Q1 Auto Aero statusCAE Overview | 9/16/2007 | © 2003 IBM Corporation

IBM Deep Computing Team

### Evolution of HPC Hardware

**MainFrames (~1979)**  
 Beginning in 1986 crash simulation drove CAE compute requirements  
 Mostly MSC.Nastran

**Vectors (~1983)**  
 SMP architecture was often first introduced in the CFD department and helped push parallel computing.

**RISC SMPs (~1994)**

**Clusters (~2001)**  
 Cluster architecture (Unix & Linux) now dominate crash & CFD environment

**Embedded (now)**  
 Embedded systems show new perspectives for CAE

5 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

### CAE System Architecture: AIX and/or Linux

- IBM business is about 50/50 Linux clusters and AIX clusters
- Linux systems tend to be special purpose
- AIX/Power systems are preferred for "implicit structures" and general purpose systems

6 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

## Challenges of Clusters

- Applications that do not scale
- Cluster nodes have weak I/O (compared to large SMP)
- Parallel I/O across the cluster
- Lots of processors generate lots of heat.
- Lots of processors generate lots of data

7 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

## Challenges of Clusters

- Applications that do not scale
  - IBM value: POWER5/6 processor and “fast” x86 MSC.Nastran
- Cluster nodes have weak I/O (compared to large SMP)
  - IBM value: MIO for Linux (fast IO libraries).
- Parallel I/O across the cluster
  - IBM value: GPFS (General Parallel File System)
- Lots of processors generate lots of heat.
  - IBM value: “blades”, “Power Executive”, water cooled rack.
- Lots of processors generate lots of data
  - IBM value: working with ISVs to promote simulation data management

8 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

**PROCESSORS**

- **Processor landscape becoming 'simplified'**
  - down to x86 and POWER
- **POWER**
  - first to dual-core
  - POWER6 to push clock
  - increasing reliance of SMT to maximize performance
  - POWER family ( i.e. embedded, gaming) influence on future
- **x86**
  - initial push to higher clocks
  - thermal problems push direction to multi-core
  - increasing reliance on SSE for HPC
  - memory future is blurry (FBdimm vs DDR?)
- **Commonality**
  - ultimately the measure of performance will be dictated by the speed and number of threads per socket.

9 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

**POWER**

- **Compute center economics paradigm shift**
  - worldwide demand for energy is increasing faster the supply
  - as a result, energy consideration with become increasingly important factor in providing CAE server solutions
  - while BG/L currently has limited applicability within CAE, its technology is pushing the envelope of energy efficiency, which will play a crucial role for future servers

2000 – Raw processing "horsepower" is the primary goal, while the infrastructure to support it is assumed ready

2006 – Raw processing "horsepower" is a given, but the infrastructure to support deployment is a limiting factor

**Three Cooling Challenges**

1. The System
2. The Rack
3. The Data Center

▪ Power and cooling spend will exceed new server spending (Gartner 2006)



Year	New server spending (US\$B)	Power and cooling (US\$B)	Installed base (M units)
1996	50	50	2
1997	55	55	3
1998	60	60	4
1999	65	65	5
2000	70	70	6
2001	75	75	7
2002	80	80	8
2003	85	85	9
2004	90	90	10
2005	95	95	11
2006	100	105	12
2007	105	115	13
2008	110	125	14
2009	115	135	15

10 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

- **Costs have become an overriding consideration in packaging**
  - product cycles changing from 36 to 6 months!
  - reuse of components is a must
- **Environmentals (space and energy) pushing packaging technology forward**
- **Utility mentality emerging**
  - shift from homogenous computer floor to constant upgrade of resource grid.

**Frames**  
↓  
**Racks**  
↓  
**Blades**  
↓  
**Bricks**

**PACKAGING**

11      2007 Q1 Auto Aero status      9/16/2007      © 2006 IBM Corporation

IBM Deep Computing Team

## **HPC Hardware Value**

*It is now much more than \$/MFLOPS*

*Total Cost of Ownership (TCO) is now more complicated.*

- **ISV application cost**
- **Power and Cooling**
- **Engineer productivity**
- **Data center floor space**

12      2007 Q1 Auto Aero status      9/16/2007      © 2006 IBM Corporation

IBM Deep Computing Team

## Contents

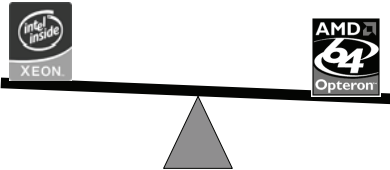
- Industry Trends
  - Auto/HPC trends
  - the 5 P's
- HPC Product Update
  - XEON 5100 series
  - Opteron "rev F"
- CAE server solutions
  - structures
  - impact analysis
  - CFD
  - Interconnects
- IBM "value add"
  - PowerExecutive, GPFS
  - MD Nastran Tuning, vMIO, Accuracy
- General Discussion

13 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

## XEON vs. Opteron Product Positioning

- The benchmark wars are in full swing
  - Intel 5160 "Woodcrest 3.0 GHz
  - AMD/Opteron "Rev. F" 2.8 GHz
  - dual-core chips with comparable performance
- It is often difficult to identify optimal product to deploy
- There are several key things to understand about each solution that help us identify which is optimal for a given workload
- But remember, the areas where there are overwhelming and compelling differences between the two (Xeon and Opteron) are usually easily identified
  - and in many cases it boils down to customer preference








14 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation



IBM Deep Computing Team

## New System x™ and BladeCenter® Servers: Intel Xeon processors




Position			
<b>x3550</b>  <i>Low cost HPC compute node</i>	<b>x3650</b>  <i>Highly available application server</i>	<b>X3850/3950</b>  <i>Mid-Market, Large Enterprise HPC</i>	<i>Ultimate scale-out integration</i> <b>HS21</b>  <i>Enterprise class scalable 2-socket blades for front and mid tier applications</i>






Key Features			
<ul style="list-style-type: none"> <li>▪ Dual socket XEON 5100 series processors</li> <li>▪ 1/32GB of FBD memory</li> <li>▪ 2(3.5") or 4(2.5") SAS internal storage</li> </ul>	<ul style="list-style-type: none"> <li>▪ Dual socket XEON 5100 series processors</li> <li>▪ 1/48GB of FBD memory</li> <li>▪ 8(2.5") or and 6(3.5") SAS + tape internal storage</li> </ul>	<ul style="list-style-type: none"> <li>▪ Four socket XEON per node</li> <li>▪ 16 2/64GB of DDR2 per node</li> <li>▪ Up to 8 nodes per system</li> <li>▪ 6(2.5") SAS internal storage per node</li> </ul>	<ul style="list-style-type: none"> <li>▪ Dual socket XEON 5100 series processors</li> <li>▪ 1/32GB of FBD memory</li> <li>▪ 2(2.5") SAS internal storage + 3(2.5") with optional SIO blade</li> </ul>

15 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

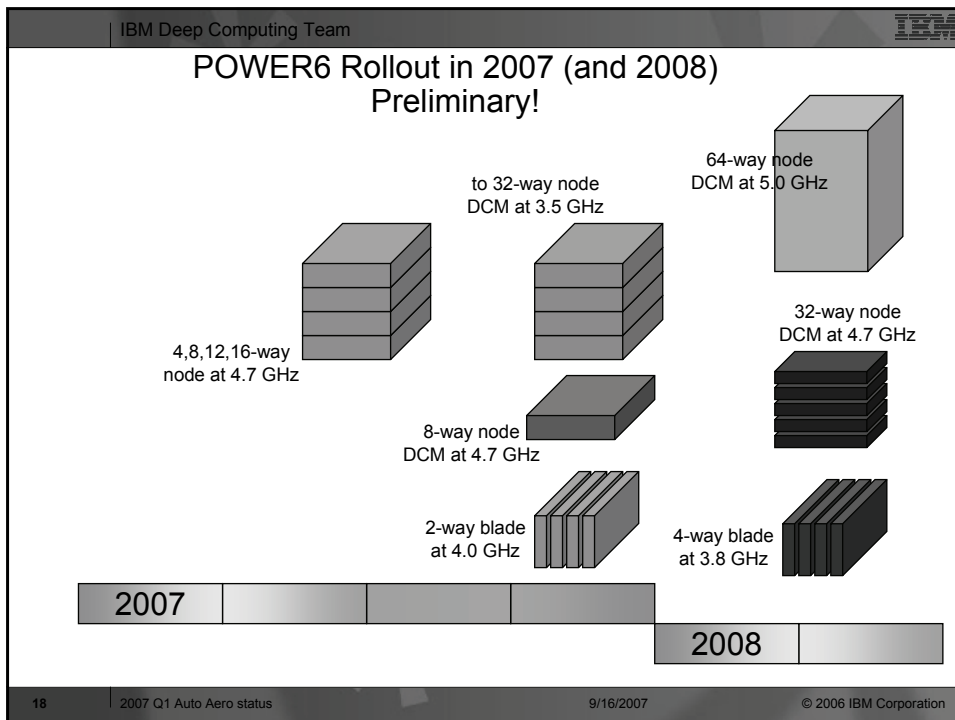
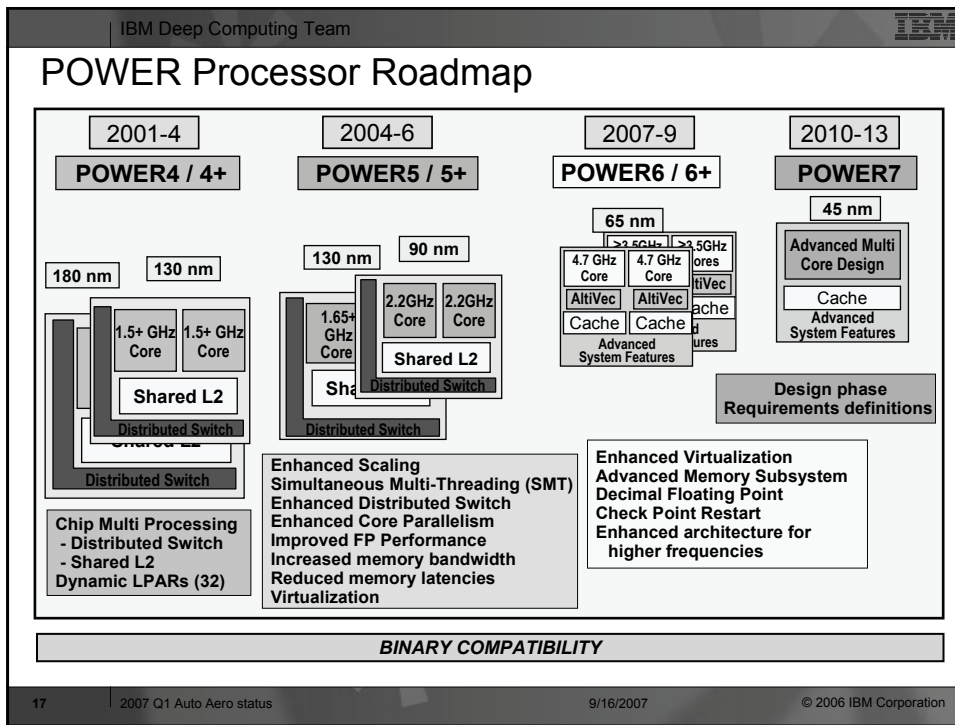
## New System x™ and BladeCenter® Servers: AMD Opteron™ Rev F



Position			
<b>x3455</b>  <i>Low cost HPC compute node</i>	<b>x3655</b>  <i>Highly available application server</i>	<b>x3755</b>  <i>Mid-Market, Large Enterprise HPC</i>	<i>Ultimate scale-out integration</i> <b>LS21</b>  <b>LS41</b>  <i>Enterprise class scalable 2-4 socket blade for front and mid tier applications</i>

Key Features			
<ul style="list-style-type: none"> <li>▪ Dual socket Opteron processors</li> <li>▪ 48GB of DDR2 memory</li> <li>▪ 3.5" Fixed SATA</li> <li>▪ Leadership I/O with PCI-E, and HTx</li> </ul>	<ul style="list-style-type: none"> <li>▪ Dual socket Opteron processors</li> <li>▪ 64GB of DDR2 memory</li> <li>▪ 2.5" and 3.5" internal storage and tape</li> <li>▪ Ready RAID and Ready RSA</li> <li>▪ Trusted Platform Module</li> <li>▪ Standard TOE</li> </ul>	<ul style="list-style-type: none"> <li>▪ Four socket Opteron processors</li> <li>▪ 128GB of DDR2 memory</li> <li>▪ 3.5" SAS internal HDD</li> <li>▪ Ready RAID and Ready RSA</li> <li>▪ Trusted Platform Module</li> <li>▪ Standard TOE</li> </ul>	<ul style="list-style-type: none"> <li>▪ Dual socket Opteron processors</li> <li>▪ 32 GB of DDR2 memory</li> <li>▪ SAS HDD technology</li> <li>▪ TOE NIC solution</li> <li>▪ High speed enablement</li> <li>▪ Supports the new SIO blade</li> </ul>

16 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation



IBM Deep Computing Team

## LSDYNA

□ Performance characteristics

- benchmarks:
  - Neon: 230K; 20ms; 4-way
  - rNeon: 550K; 30ms; 4-way
  - 3Car: 1.2M; 10ms; 4-way
- standard benchmarks: [www.topcrunch.org](http://www.topcrunch.org)
- cache friendly (follows SPECfp)
  - similar performance on POWER5, Opteron and XEON
  - JS21 offers potential for excellent price/performance for AIX customers
- scales well with clusters
  - high performance network now the norm

□ IBM solutions

- when AIX is the most important factor
  - System p JS21; 2.5GHz 4-way nodes; 8GBmem; 1 internal drive; Myrinet preferred
  - System p p5 575+; 1.9GHz 16-way nodes; 16GBmem; 2 internal drives; HPS
- when price/performance is most important factor
  - System x HS21; 3.00 GHz 4-way blades; 8GBmem, 1 internal drive; Myrinet or IB
  - System x 3550; 3.00 GHz 4-way nodes; 8GBmem, 1-2 internal drives; Myrinet or IB

Performance (relative to x336, higher is better)

Benchmark	x336/3.6	WC/3.0	Is20/2.2	js21/2.5
Neon	1.0	1.5	1.1	1.0
rNeon	1.0	1.5	1.2	0.8
3Car	1.0	1.4	1.1	0.9

x336/3.6
  WC/3.0
  Is20/2.2
  js21/2.5

**LSTC**  
Livermore Software Technology Corp.

19
2007 Q1 Auto Aero status
9/16/2007
© 2006 IBM Corporation

IBM Deep Computing Team

## LS-DYNA comparison: dual-core vs. quad-core

Sockets	dual-core 3.0 GHz	quad-core 2.66 GHz
2 sockets	~41000	~34000
4 sockets	~21000	~18000
8 sockets	~11000	-


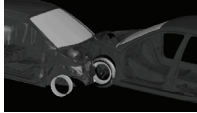
April, 2007 testing, 3-car model, 795k elements, 150 msec  
 IBM x3550 3.0 GHz Xeon 5160 "Woodcrest"  
 IBM x3550 2.66 GHz Xeon X5355 "Clovertown"

20
2007 Q1 Auto Aero status
9/16/2007
© 2006 IBM Corporation

IBM Deep Computing Team

## CAE Server Solutions

- No dominant server choice for all CAE applications
  - System p
    - strength of AIX
    - industry leading performance for many problems
    - well balanced performance for wide variety of simulation
  - System x
    - economics and flexibilities of open standards
    - extensive application portfolio
    - typically excellent price/performance
- No dominant server strategy for CAE customers
  - General purpose CAE servers
    - System p typically offers best performance for variety of applications
    - System x typically offer best price/performance
  - Application specific CAE servers
    - complex landscape which is always in flux

21 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

## Contents

- Industry Trends
  - Auto/HPC trends
- HPC Product Update
  - XEON 5100 series
  - Opteron "rev F"
- IBM "value add"
  - GPFS
  - vMIO
- General Discussion

22 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

## Scalable Parallel I/O: General Parallel File System (GPFS)

- **NFS**
  - Client-server file systems have server bottleneck and protocol overhead
- **SAN**
  - SAN with a single metadata server have potential bottleneck
- **GPFS**
  - **General Purpose**
    - Any node can read from or write to any of the disks
    - The entire cluster can be administered from a single node
    - Supports Linux, AIX and mixed clusters
  - **High Performance**
    - Has provided 15GB/s to a single node and 100GB/s against a single file
    - *GPFS is not a client-server file system and has much lower protocol overhead*
    - All system data & metadata is equally accessible from all nodes
    - All data & metadata flows between the disks and nodes in parallel
  - **Scalability**
    - Currently supports 100s of nodes and 200+TB of storage over LAN or HPS (more by special bid)
  - **Reliability**
    - Parallel operation means no single point of failure
    - One large research customer reported 100% uptime for GPFS for an entire year

23 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

## Members of IBM's CAE Team (1)

<u>Nick Alsopp</u>	ABAQUS, CFX	Many years experience in HPC and with HKS
<u>Balaji Atyam</u>	ANSYS, LMS, Madymo	Application Support USA
<u>John Bauer</u>	HPC I/O Libraries	Original developer of EIEIO libraries, 19 years HPC experience
<u>Steve Behling</u>	STAR-CD	13 years experience with STAR-CD source code
<u>Achim Bömelburg</u>	Permas, CAE Team	16 years experience with automotive customers
<u>David Wei Chen</u>	Detroit CAE Team	9 years working with automotive users in Detroit
<u>Greg Clifford</u>	Leader CAE Practice	20 years working with CAE customers and ISVs
<u>Martin Feyereisen</u>	Pam-Crash, LS-Dyna	Considered a member of the ESI development team
<u>John Hague</u>	VECTIS	ACTC Team UK

24 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

## Members of IBM's CAE Team (2)



<u>Holger Holthoff</u>	RADIOSS, AVL	11 years experience in parallel computing on IBM platforms
<u>Nobuhiko Kudanami</u>	Tokyo CAE Team	Working with automotive users in Tokyo
<u>Guangye Li</u>	LS-DYNA	Extensive experience with LS-DYNA on Linux
<u>Doug Petesch</u>	NASTRAN, AMLS	15 years experience working with MSC and customers
<u>Hari Reddy</u>	FLUENT, PowerFLOW	6 years experience with FLUENT, experience with various CAE codes
<u>A. Sugavanam</u>	PowerFLOW, CEM	Many years experience with NASA and CAE codes
<u>Erling Weibust</u>	CAE Team Sweden	19 years working with technical customers in Nordic region
<u>Jeff Zais</u>	Leader CAE technical team	Key player in the 1999 success of MPP-DYNA
...		

25 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation

IBM Deep Computing Team

## The IBM Value

- **Experienced HPC applications team**
- **Worldwide customers**
  - Longstanding relations with key application vendors
- **Full range of computing solutions**
  - POWER6 to Linux Clusters
  - Storage Solutions
- **Presence of IBM**
  - Stable company, growing in technical computing
  - Able to offer complementary solutions for storage and the desktop.

26 | 2007 Q1 Auto Aero status | 9/16/2007 | © 2006 IBM Corporation