HPC Clusters and Solutions

Alberto Galli HPC Consultant - Presales Italy Mobile: +39 335 6322966 alberto.galli@hp.com September 2010



©2009 HP

CHALLENGES TO INNOVATION

CUSTOMER OBJECTIVES





OBSTACLES

- Increasing scale and complexity of deployments and applications
- Exploding volumes of data
- Pressure to provide instantaneous results
- Affordability
- Power and space constraints
- Technology and market risks and uncertainty
- Adapting to new business models and technologies



TODAY'S REQUIREMENTS FOR THE HPC DATA CENTER

Latest technology (e.g., multi-core, QDR IB)

Compute and data scalability

Energy efficiency

Flexible deployment models for the data center

Data center and HPC expertise

Enabling Affordable Breakthroug h Innotation



HP UNIFIED CLUSTER PORTFOLIO

HPC Technical and Enterprise	
------------------------------	--

HPC application, development and cloud software portfolio

Advanced	and	specialty	options

(Accelerators, Visualization, other)

Sool	a h		mon	00	om	ont
JUd				au	еш	еш
				~ 3	••••	· · · ·

(HP x9000 NSS, Lustre Cluster FS)

Cluster management layer

Partner and Open Source choice Microsoft Windows HPC Server 2008

Operating environment and OS extensions

Linux

HP CMU

Windows

HP cluster platforms

HP ProLiant servers, HP BladeSystem, multiple interconnects

HP Datacenter Products & Services



HPC CLUSTER MANAGEMENT OPTIONS FOR CLUSTER PLATFORMS



Available factory-installed with Operating System and HP-MPI

DRAMATIC PERFORMANCE BOOST WITH PROLIANT G6 SERVERS



Performances are measured on specific HP server configurations and specific data sets . Any difference in system configurations or data sets may affect actual performance.



OPERATION SAVINGS WITH PROLIANT G6 SERVERS

Power Saving** over G5 for the same amount of work



** Power savings are estimated based on BladeSystem Power Calculator.

* Other names and brands may be claimed as the property of others.

7 ©2009

INTERCONNECT SPECTRUM FOR SCALE-OUT SOLUTION

- Gigabit Ethernet dominates in total volume
 - Pervasive mature technology
 - Most cost effective for workloads not interconnect latency and bandwidth bound
- 10GE being increasingly deployed for scale-out solutions
 - Decreasing cost: 10G LOM, switch cost, etc
 - Performance improvements and acceleration technologies (TOE, RNIC, RoEE?)
 - Most for higher bandwidth requirement of network traffics
- InfiniBand remains the choice for HPC clusters
 - Extremely low latency, high bandwidth, scale to 1000s nodes
 - Most for MPI-based applications, low latency market data systems, scalable data warehouse applications



HP PROLIANT SERVERS OPTIMIZED FOR EXTREME SCALE

-DL1000

- Designed for scale-out from SMB to extreme scale-out
- Up to 4 nodes in a high efficiency 2U chassis
- Non-plugable servers, with front hot-plug hard drives
- Standard racks, with traditional rear cabling



HP ProLiant SL1000 Scalable System

-SL6000

- Designed for extreme scale-out deployments by the rack
- High efficiency, modular power and cooling chassis
- Front serviceable server/storage trays de-coupled from the infrastructure
- Standard racks, with front I/O cabling



Highly Flexible s6500 Chassis

Multi-node, Shared Power & Cooling Architecture



- Shared Power & Fans
- Optional Hot-Plug Redundant PSU
- Energy efficient Hot Plug fans
- Single or 3 Phase Load Balancing
- 94% Platinum Common Slot Power Supplies
- N +1 Capable Power Supplies (up to 4)

Benefits: Low cost, high efficiency chassis

- 4U Chassis for deployment flexibility
- Standard 19" racks, with front I/O cabling
- Unrestricted airflow (no mid-plane or I/O connectors)
- Reduced weight
- Individually Serviceable Nodes
- Variety of optimized Node Modules
 1200mm rack to close doors
 - SL Advanced Power Manager Support
 - Power Monitoring
 - Node Level Power Off/On





HP ProLiant SL390s Fast Fabric Compute

Tawighth 1U tall, 2 hot plug nodes in 1U



囫

HP ProLiant SL390s, 0 to 3 GPUs

¹/₂ width, 2U tall (effective density=1U); 4 trays per 4U chassis



WHY LOOK AT ACCELERATORS?

- Facing computing challenges
- Need to innovate to meet computing needs
- Need to increase performance
 - While improving power, cooling, and space usage
 - While improving price/performance
- Accelerators offer an alternative to multi-core for meeting performance needs
- Accelerators add heterogeneous co-processors to standard systems to improve performance
 - More silicon doing what you need



PROS AND CONS OF HPC ACCELERATORS Pros

- Hundreds of functional units executing in parallel
- Speedup applications by 2x, 10x, 30x, 100x or more!
- Really fast
- Excellent for a growing number of highly parallel applications
 When it works, it can really fly!

Cons

- Can be hard to program
 - But getting easier to program every quarter
- Not useful for many general purpose applications



HP SCI ACCELERATOR PROGRAM

- HW Accelerators provide substantial speed-up for certain applications
 - GPUs, FPGAs, custom processors
- HP Accelerator Program
 - partner with accelerator vendors
 - enable accelerators in ProLiant platforms
 PCle x16 slots, power, BIOS
 - qualify selected accelerators
 - Nvidia and AMD/ATI GPUs, Nallatech & XDI FGPAs, ClearSpeed and Mercury Cell custom processors
 - deliver integrated solutions
 - Proliant DL160seG6 (1U 2P Xeon) + Nvidia Tesla S1070 (1U 4 GPUs)
 - Configure DL160seG6 with 2, 4, or 6 GPUs







WHERE WE SEE ACCELERATORS





MULTI-CORE AND APPLICATIONS: THE OPPORTUNITY AND THE RISK

- Opportunity more raw compute power
 - Multi-core processors can deliver significant performance benefits and relieve compute cluster sprawl.
 - Each core contains its own set of execution resources
 - Results in very low latency parallel execution of application threads within a single physical CPU package.
- The risk it's a lot more complicated
 - An example:
 - It's possible for one memory-intensive job to saturate the shared memory bus resulting in degraded performance for the system.
 - As the number of cores per processor and the number of threaded applications increase, the performance of more and more applications will be limited by the processor's memory bandwidth, I/O etc



HP MULTI-CORE OPTIMIZATION PROGRAM



- It is a vehicle to identify, scope and then address the issues that Industry Standard Multi-core Systems present to our customers, partners and the industry.
- Best of breed program drawing resources from:
 - Across HP including HP labs
 - Our Technology partners The intent of the Multi-core program is to
- Increased visibility and effective usage of Multi-core high performance systems
- Maximizing of total application/job throughput in multi-core systems
- Maximizing application performance in Multicore systems, using
 - Traditional parallelization and multi thread techniques for compilers and debuggers
 - New and emerging technologies
 - Enhancing legacy HPC math and libraries for mutli-core



HP X9000 ADDRESSES:

Scalable tiers of storage Example: Health and Life Sciences	 Problem Massive scalability required at a moments notice Data availability during system expansions Performance degradation over time Infrequently accessed data on high performance disks
Parallel workloads Example: Media & Entertainment	 Problem Lots of small files require concurrent access Limited bandwidth into storage system Controller hardware bottlenecks Performance - manageability tradeoff
Disaster recovery Example: Financial Services	 Problem Copying very large files concurrently to a remote site Cost and complexity System performance impacts Unused storage at remote location



COMPLEMENTARY SCALABLE STORAGE SOLUTIONS FOR HPC

X9000 Network Storage System and Fusion File System

- Shared datacenter multipurpose storage
 - -Desktops, clusters, farms, clouds
 - -Windows, Linux; CIFS/NFS
- High performance and scale with distributed metadata
- Data tiering
- Continuous replication
- Popular for web 2.0; bioInformatices, FSI, GeoScience apps with many small files



Cluster File System with DirectData Networks

- Tightly coupled HPC storage
 - For large HPC Linux clusters with large single files/single stream requirements (traditional HPC, such as CAE)
- High parallel bandwidth
- High capacity
- High scalability and reliability
- Lustre-open source technology



IBRIX Fusion Glossary

Segment: Storage Element

- An atomic storage element in the IBRIX File System
- A logical storage volume
- A file and directory storage bucket
- Segments can be added dynamically

Segment Servers: File Serving to clients, compute nodes and app servers

- Manage elements for specific segments
 - Data allocation and layout
 - Metadata
 - Locks
- Provide file system NAS functions
 - Protocols: NFS, CIFS, IBRIX Client
 - File System and/or Segment Cache
- Provides file system wide coherent cache
- Communicates with other segment servers
 - For metadata information/operations
 - For file information and access (NFS & CIFS)
 - For cluster status and other information



Fusion Manager: Management

- Cluster management system
- Web based and CLI interface
- Allocates segments to servers
- Manages faults & segment servers
- Provides logging and notification

SCALABLE VISUALIZATION FOR HPC

Remote visualization

Remote access to graphic enabled servers over a standard network Simplify management and share resources

Scale-up beyond the workstation

Process large models Drive display to multi-tile panels and caves

Simplify HPC visualization

VizStack manages surFace displays

Supports launch of parallel visualization applications







Thank you

