

## Blades—An Emerging System Design Model for Economic Delivery of High Performance Computing

. LUN %UHVQLNHU  
6\WHPV \$UFKLWFW  
,QWUQHWDQG \$SSQLFDWLRQV 6\WHPV /DERUDWRU\  
+HZOHWW3DFNDUG

6HSHPEHU



## Agenda

### Agenda

,QWURGXFILRQ ' HILQILRQV  
 ZULJLQV RI WKH %ODGH 6HUYHU  
 /LPILQILRQV RI )LUVW \* HQHUDILRQ ' HVLJQV  
 6HFRQG \* HQHUDILRQ ' HVLJQ  
 )XIXUH ' LUHFILRQV



During the recent rise and subsequent fall of the internet bubble, a new computer system design model emerged, primarily from venture capital start-ups. Bladed Systems, dense arrays of single board computers housed in a common chassis, seemed a promising way for service providers to keep pace with the anticipated dot.com inspired build-out. The blades were dense, low bandwidth and low in computational power, but they were suited to rapid deployment of massive content delivery.

Along with other lessons learned as the irrational exuberance faded and unviable business models and their edge applications were winnowed from data centers, the designers of bladed systems began to realize that blades had the potential to move from edge-only applications into high performance enterprise, communication, and technical compute.

This discussion will cover the origins of the bladed server system model, the design, and limitations of the first generation designs, the ongoing technology changes which are enablers to second generation blades, their application to a wider set of multi-system compute problems, and some possible directions for future development.

## A Blade is ...

### Definitions

\$Q LQFØXVLYH FRPSXILQJ V\VIHP WKDW LQFØXGHV SURFHVVRI  
PHPRU\ QHWZRUN FRQQHFILRQV DQG DWRFDWHG HØFWUROLFV RQ D VLOJØH  
PRIKHUERDUG

7KH VHUYHU EODGH \NSLFDØ\ LV DWRFDWHG ZLWK DQ HOFØRVXUH V\VIHP WKDW  
DØRZV PXOWSOH EODGHV IR EH KRXXHG LQ D VIDQGDUG VHUYHU \VXE UDFN• RU  
HOFØRVXUH WKDW VKDUH UHVRXUFHV VXFK DV SRZHU VXSSØHV DQG FRRQLQJ  
IDQV

%ODGHV «

‡ DUH HDVLØ\ DFFHWLEØH

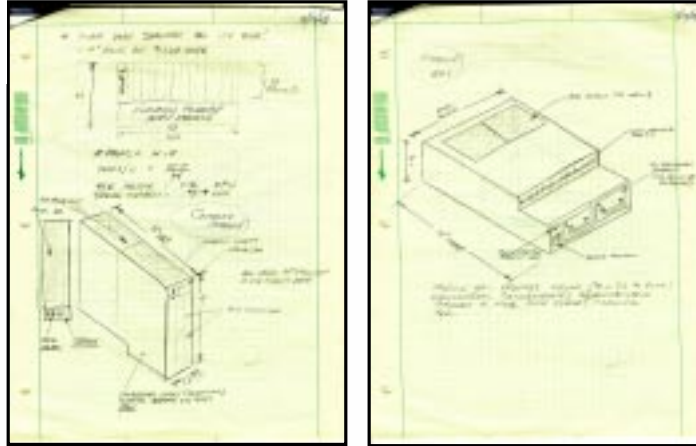
‡ RIIHU LQFUHDVHG FRPSXILQJ GHQVLØ\

‡ KDYH PRGXØDU DUFKLWHFWXUH WKDW HQVXUHV IØH[LELØ\ DQG VFDØDELØ\

Source: Provisioning the Internet Infrastructure: Server Blades and Dynamic Workload Management IDC  
Document # 24155



**First, some definitions. According to IDC, A Blade is ...**



It is, perhaps, instructive to dig back to the origins of the systems packaging concept that has come to fit this definition as it developed as a concept within Hewlett-Packard.

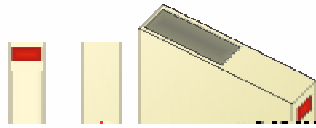
In late 1998, I went on a customer visit to the primary data center of a national internet service provider in Fairfax County, Virginia.

I had just completed architecting the first 2U HP-UX PA-RISC rack optimized server, after having lead a team planning for the initial turn on of what eventually became HP's high end SMP SuperDome systems.

Although the new 2U server was designed to be a general purpose system, the nascent Internet Service Provider market was also a target. In designing a dense system, there were tradeoffs made to accommodate the shrinking form factor. I was intrigued by what might be accomplished if we left the general purpose server requirements behind, concentrated solely on the emerging ISP environment and marketplace, and their needs for networked computing. I tried to imagine the minimal system to bring packets in a narrow pipe, run them through an application and spit them out a fat pipe. I also wondered if the high end SMP mesh backplane technologies I had been exposed to might be migrated down cost effectively.

On the flight out I had a chance to sketch up some ideas I had been mulling over for several months, and came up with a system packaging concept similar to that we would today call a blade server. My visit to the ISP data center, especially contrasting it against my experiences in traditional Enterprise data centers, reinforced my belief that this alternative packaging concept held promise.

## Initial Blade Sketches

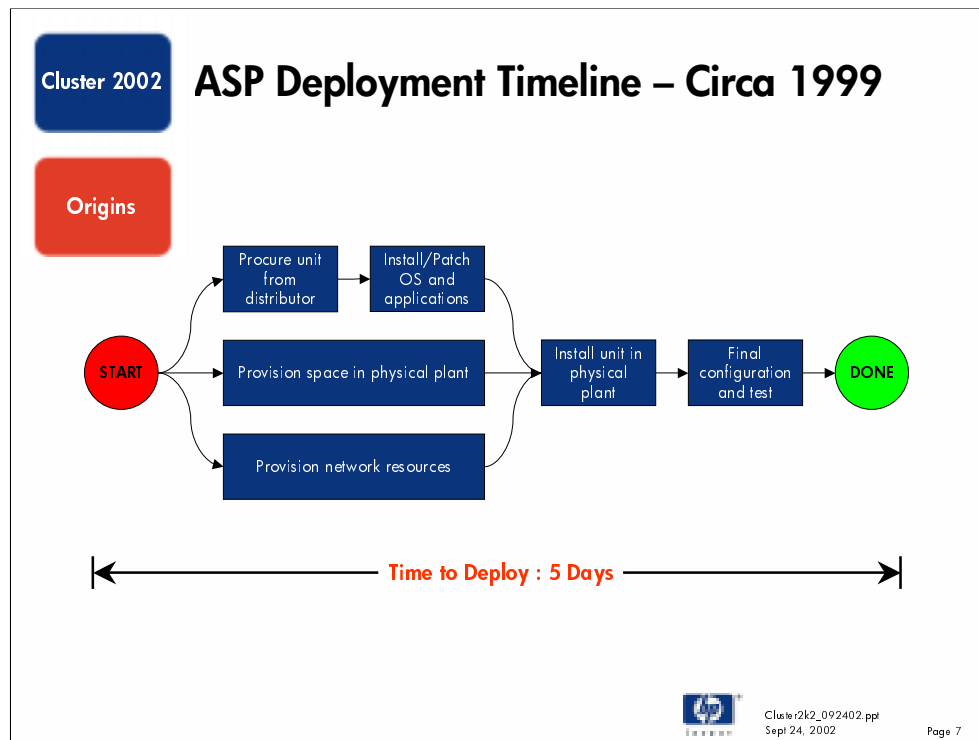


## Dominant xSP Concerns – Circa 1999

**Time-to-Deployment** &RPSHWLWRO LV FXWIKURDW +RZ  
TXLFN\ FDQ D QHZ VHUYLFH EH GHSUR\HG LQ WKH SK\VLFD SODQ'

**System Density** 7KH GDWD FHOWHU LQFXUV FRVW LQ DUHD DQG  
UHFRYHU FRVW LQ V\WHPV DUHD +RZ PDQ\ V\WHPV FDQ EH





For example, at that time, a cutting edge ASP was able to deliver a new hosted application server in five days and were looking to shave off half hours from that time.

## Dominant xSP Concerns – Circa 1999

### Origins

**Time-to-Deployment** &RPSHWLWLRQ LV FXWIKURDW +RZ  
TXLFN\ FDQ D QHZ VHUYLFH EH GHSUR\HG LQ WKH SK\VLFD SODQ'

**System Density** 7KH GDWD FHQWHU LQFXUV FRVW LQ DUHD DQG  
UHFYHUV FRVW LQ \V\WHPV DUHD +RZ PDQ\ \V\WHPV FDQ EH  
EXUGHQHG RQ HDFK VTXDUH IRRW RI GDWD FHQWHU"

**Moving Customers up the price list** &XUHQW\ VKDUHG  
KRVLQJ RQ\ VHOO DW WKH LQJRGXFWRU\ VHUYLFH OHYHO ORYLQJ  
FXVRPHUV XS WKH SULFH QLVW UHTXLUHV GHGLFDWHG KDUGZDUH



The second was that the xSP data center incurred costs per square foot of managed data center and recovered revenue in systems per square foot. The final common theme was that virtual hosting on a shared system was already a lost leader, and real revenue required dedicated hardware.



## First Generation Blade



(DFK E0DGH FRQNDLQV

‡ O+] 3,,, . /9

‡ O% \*% 3&  
6' 5\$O

‡ \*% ORELOH , ' ( ' ULYH

‡ ' XD0 1,&V

‡ ,QWHJUDWHG ODQDJHPHQW

(DFK 8 FKDWLV FRQNDLQV

‡ 6HUYHU E0DGHV

‡ 1 3RZHU

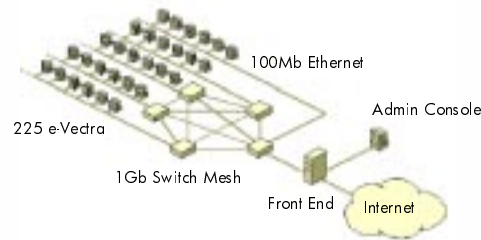
‡ ' XD0 (WKHUQHW 6ZLWFKHV



Cluster 2002

## Mainstream Technology Clusters

First  
Generation



&0XVWU 7RSR0RJ\

ODLQVWHD P 7HFKQR0RJ\ &0XVWU DFKLHYHG UDQN RI WK RQ 723 DV RI



Cluster2k2\_092402.ppt  
Sept 24, 2002

Page 10

As a quick aside, It should be noted that at even this level of performance, bladed systems could be a substantial asset to the creation of cost effective, easily maintainable clusters. In October 2000, Richard Bruno, a research at HP Labs in Grenoble, who will be presenting later today, teamed with INRIA Rhone-Alpes to construct a cluster using 225 HP e-Vetra business desktop PCs and HP Procurve ethernet switches and reached the 385th rank in the TOP500 list by June 2001. The individual e-Vetra memory, CPU, and network configuration was very similar to the first generation of blades. But, imagine how much more attractive a single rack of blade servers would be to a data center manager than the open frame rack and hand wiring involved in the setup of a 'mainstream technology' cluster.

```
[root@imageserver]# getimage -g my-golden-client -image web_server_image_v1
This program will get the "web_server_image_v1" system image
from "my-golden-client"
making the assumption that all filesystems considered part
of the system image are using ext2, ext3, or reiserfs.
This program will not get /proc, NFS, or other filesystems
not mentioned above.
See "getimage -help" for command line options.
Continue? ([y]/n): y
Retrieving /etc/systemimager/mounted_filesystems from mygolden-
client to check for mounted filesystems...
----- my-goldenclient
mounted_filesystems RETRIEVAL PROGRESS -----
receiving file list ... done
/var/spool/systemimager/images/web_server_image_v1/etc/systemimager/mounted_wrot
e 112 bytes read 294 bytes 852.00 bytes/sec
total size is 180 speedup is 0.42
----- my-goldenclient
mounted_filesystems RETRIEVAL FINISHED -----
Retrieving image web_server_image_v1 from my-golden-client
----- web_server_image_v1 IMAGE RETRIEVAL PROGRESS --
receiving file list ... done
```



Offering maximum systems/square foot of dedicated hardware, the blade systems seemed poised to maximize the revenue return per square foot of a network edge xSP data center. And, in order to minimize the time-to-deployment, most blade system providers were also beginning to offer provisioning software to allow for simplified mastering of the on-blade OS images via a central management station. For low end, low bandwidth applications, blades now had traction in every part of the time-to-deployment problem that was first evident in 1999.

However, there were some inhibitors. At the time, single chip 10/100 ethernet switching technology was only just becoming commodity. As a result, some first generation designs lack an integrated switch. For these systems the promise of easy deployment and maintenance of systems was true at the blade level, but not at the chassis level.

Second, the storage subsystems included in the first generation blades based on mobile style drives were more expensive, lower performance, and lower reliability than their general purpose server counterparts. The prevalent 10/100 networking didn't really have the bandwidth to allow for network attached storage, and no first generation blades had storage area network interfaces. Finally, the blades were limited in overall CPU performance and memory array size due to power and space considerations.

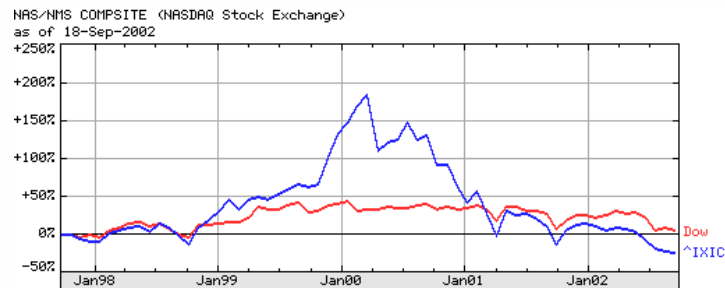
Another problem was that some designs provided little more for manageability than a command line or GUI front end to manage the images for LAN boot protocols. For those designs without a solid, blade-centric manageability model, the nearly 6X increase in density only exacerbated Server Sprawl, the proliferation and deployment of massive numbers of systems that can lead to acute manageability problems and eventually to unsustainable deployments.

Still, as long as web sites need to be provisioned as quickly as before ...

Cluster 2002

## What New Economy?

Limitations



Cluster2k2\_092402.ppt  
Sept 24, 2002

Page 12

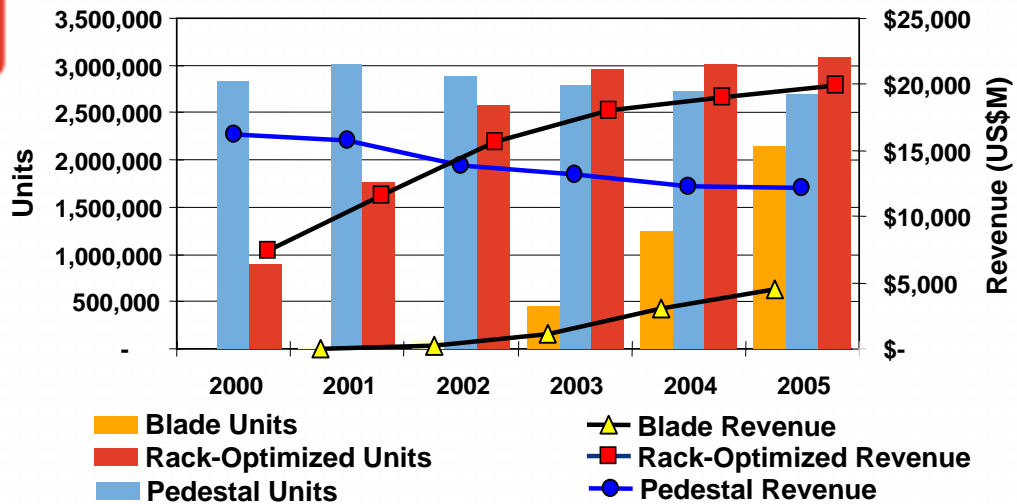
**Pop! Did we say new economy? With the very rapid decay of the internet bubble economy, the demand for rapid build out of thin, dense servers was quickly eroded. Along with failing dot.coms and xSPs, many of the startups designing bladed systems also found themselves without a market.**

**But, this shakeout did not dissuade the major vendors from continuing to explore the density and total cost of ownership advantages of blades. In fact, although thin blades are a great way to deploy network edge servers, the cost of ownership advantages of blades can be a benefit in high performance systems as well.**

Cluster 2002

## Bladed Architectures: What's the Opportunity?

Second  
Generation



By 2005, IDC believes the blade form factor will capture approximately 27% of entry server unit sales and 12% of entry server revenue.

Source: IDC August 2001



Cluster2k2\_092402.ppt  
Sept 24, 2002

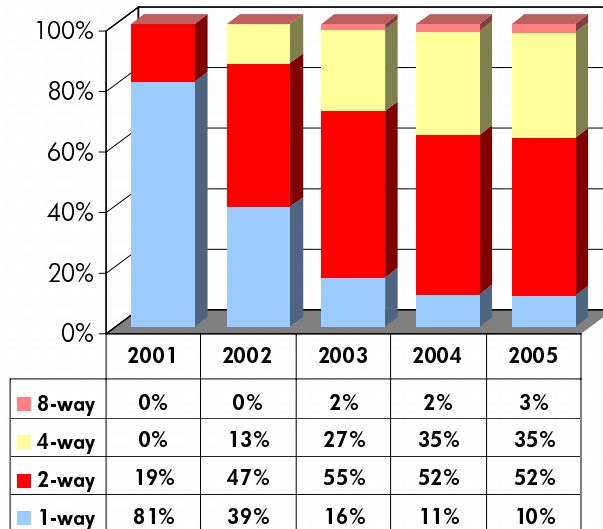
Page 13

As is shown in this IDC market projection, the penetration of blade into the server marketplace is anticipated to ramp steadily over the next several years, culminating in a 27% penetration rate in 2005.

Cluster 2002

Second  
Generation

## Bladed Architectures: Where Are the Sweet Spots?



EL00LRQ

Share of bladed server revenue quickly migrates to the 2 & 4 processor  
"flavors" as volume ramps and commoditization occurs

Source: IDC August 2001



Cluster2k2\_092402.ppt  
Sept 24, 2002

Page 14

But, perhaps also more interesting is this IDC product segmentation projection which shows the penetration of blades with dual, quad, and even octal SMP configurations in the same timeframe. Clearly these will not be the mobile technology enabled blades of the first generation, but what will they look like?

## Second Generation Blade

Second  
Generation



- ‡ 3URFHVVRIJ ORELOH /9 8/9 → ODLOLOH VHUYHU 3URFHVVRIJ
- ‡ 6O3 6LOJOH → ' XD0 RU PRUH 3URFHVVRIJ
- ‡ OHPRU\ VLJHV 6XE \*% SHU SURFHVVRIJ → \* UHDIHU WKDQ \*% SHU SURFHVVRIJ
- ‡ (WKHUQH 1 ,&V ' XD0 → OXOLSOH
- ‡ ODQDJHDELQW UDFN RSILPLJHG VHUYHU OHYHJDJHG → EODGH DZDUH
- ‡ /RFD0 6VRUDJH ORZ SHUIRUPDQFH PRELOH → KLJK SHUIRUPDQFH KLJK UHDELQW VHUYHU VRUDJH
- ‡ 5HPRUH 6VRUDJH /RZ SHUIRUPDQFH 1 \$6 → KLJK SHUIRUPDQFH 1 \$6 RU 6\$1
- ‡ &KDWLV IRUP IDFWRU 8 → 8
- ‡ %ODGHV IRUP IDFWRU LOFUHDVH LQ GHSWK ZLGWK DQG KHLJKW



This example of an HP ProLiant BL p-class blade is generally indicative of what most major vendors are currently delivering or planning for this next generation of high performance blades. The general theme of these designs is that the same technology present in high performance 1U/2U rack optimized servers will be available in this form factor. In particular:

- ‡ ORELOH /9 8/9 SURFHVVRIJ PLJUDIH WR PDLOLOH VHUYHU SURFHVVRIJ
- ‡ 6LOJOH SURFHVVRIJ PLJUDIH WR ' XD0 RU PRUH 3URFHVVRIJ
- ‡ OHPRU\ VLJHV PLJUDIH IURP 6XE \*% SHU SURFHVVRIJ WR JUHDIHU WKDQ \*% SHU SURFHVVRIJ
- ‡ (WKHUQH 1 ,&V PLJUDIH IURP ' XD0 WR OXOLSOH
- ‡ ODQDJHDELQW PLJUDIHV IURP UDFN RSILPLJHG VHUYHU OHYHJDJHG WR EODGH VSHFLILF
- ‡ /RFD0 6VRUDJH PLJUDIHV IURP ORZ SHUIRUPDQFH PRELOH WR KLJK SHUIRUPDQFH KLJK UHDELQW VHUYHU VRUDJH
- ‡ 5HPRUH 6VRUDJH LV HQDEOHG YLD KLJK SHUIRUPDQFH 1 \$6 RU 6\$1
- ‡ &KDWLV IRUP IDFWRU LOFUHDVH IURP 8 WR 8
- ‡ %ODGHV IRUP IDFWRU LOFUHDVH LQ GHSWK ZLGWK DQG KHLJKW

This provides a snapshot of currently available high performance bladed systems. It is interesting that although the first generation system designs responded to the economic drivers of the internet bubble, these mainstream performance systems, through their vastly decreased MTTR, their excellent blade-centric management integration, and their amortization of shared resources over multiple systems, will drive blade adoption past the network edge and into the enterprise data center. As the bladed systems model matures, and especially as innovation continues in the management software infrastructure, modular hardware will continue to penetrate traditional rack optimized server market. I say modular hardware, since the label blade will become an increasingly inaccurate metaphor for future higher performance modules

Cluster 2002

## Rapid Deployment Tools

Second Generation

The screenshot shows the HP Proliant Essentials Rapid Deployment Pack interface. It features a tree view on the left for organizing servers, a central pane for running events, and a right pane for event details. Red callout boxes highlight key features:

- Flexible organization and grouping of all servers:** Points to the tree view on the left.
- View physical location of blade servers:** Points to a specific server entry in the tree view.
- Run an event by simply dragging it to a server or group (or drag servers or groups to the event):** Points to an event in the central pane.
- Use the Event Wizard to easily set-up tasks:** Points to the 'Event Wizard' button in the top right.
- Run scripts and programs on remote servers, use remote console, copy files, control power:** Points to the 'Details' section of an event.
- See details about your servers, groups, and events / tasks:** Points to the 'Details' section of an event.

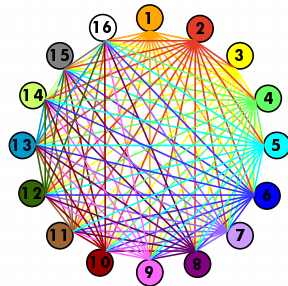
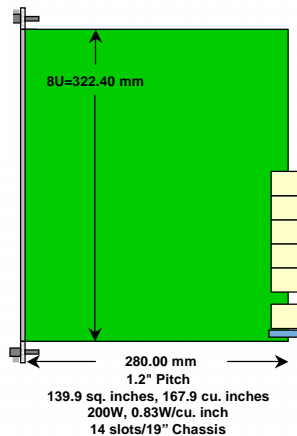
Cluster2k2\_092402.ppt  
Sept 24, 2002  
Page 16

Along with Second generation blades will come second generation rapid deployment tools. This screenshot of HP's Proliant Essential Rapid Deployment Pack shows the advances in blade-specific management tools that can prevent Server Sprawl.



# Advanced Telecom Computing Architecture

Futures



(DFK E0DGH FROWDLQV

‡ \$ : 7KHUPD0 (QYHORSH HQDE0HV ' XD0  
VORW ' 3 ,3) RU 43 , \$

‡ )X00\ FRQQHFVHG OHVK ' DWD 7UDQVSRUW  
3URWRFR0 \$JQRVLF )DEULF (DFK /LQ FDQ  
VXSSRUW IRU H[DPS0H

† ; ,QILOLEDQG 0LQN ; ,QILOLEDQG 0LQV

† %; 0LQV

† \* E( ; \$8, 0LQN

† \* SEV )& 0LQV

‡ ' HGLFDVHG ODQDJHPHQW ' XD0 6INDU  
(WKHUQHW

‡ SDU \$GMDFHQW 60RW 8SGDWH EXV

‡ 5HGXQGDQW E0DGH FHQWLF KDUGZDUH  
PDQDJHPHQW EXV

‡ 1 5HGXQGDQW 3RZHU

‡ &RVW IRFXVHG PRGX0H PHFKDQLFD0 GHVLJQ



Cluster2k2\_0924.02.ppt  
Sept 24, 2002

Page 17

So blades systems are now capable of bringing their total cost of ownership and rapid deployment characteristics to the dense, horizontally scaled network edge as well as the vertically scaled data center. Where might they go next?

A possible leading indicator of trends in modular hardware features is an effort that I have been driving for HP in the last year, the Advanced Telecommunication Computing Architecture or ATCA draft standard. ATCA is a draft standard of the PCI Industrial Computer Manufacturers Group targeted at bringing next generation modular design standards to a converged communications compute market and is slated for adoption by year end. This is a summary of some features of the current draft definition of ATCA. While some aspects of the draft standard are particularly optimized to the converged communications market, many of the features of the standard could be leveraged into a third generation blade.

# ATCA Target Serial I/O Standard Comparisons

## Futures

Serial Standard	Data rate per Channel (Gbps)	Baud rate per Channel (Gbps)	Ref Clock (MHz)	Pairs/Channel Channel Configs
Infiniband	2.0	2.5	125	2 1X, 2X, 4X, 12X
1Gb Ethernet 1000BASE-CX	1.0	1.25	62.5	2 1X
10Gb Ethernet XAUI	2.5	3.125	156.25	2 4X
Fibre Channel	0.85 / 1.7	1.06 / 2.12	53.125 / 106.25	2 1X
Serial ATA	1.2	1.5	75	2 1X
Serial RapidIO	2.5	3.125	156.25	2 1X, 4X
PCI Express	2.0	2.5	125	2 1X-32X



One interesting aspect of ATCA is that while it is optimized around modular hardware interconnected via a fully connected mesh of high speed serial I/O connects, it doesn't specify which I/O connects are utilized. The integrated managements systems handles the possible heterogeneous mixture of technologies. This chart shows a current snapshot of ATCA compatible I/O connects. All these I/O connects have very similar LVDS signaling, similar bit counts and bit rates, and differ primarily in the lower layer protocol enabling and end use models.

Now this brings up an interesting question: what will the fabric of choice be for the next generation bladed system. Ethernet was the obvious choice for first generation systems, since their market called for and they only had the capacity for network connectivity. The second generation is still predominately ethernet based, although now with the capacity for high performance storage interconnect, fiber channel has been discussed. Also, from Infiniband silicon vendors have demonstrated bladed systems.

At this point, it is difficult to deny the economy of scale being afforded to ethernet is a distinct advantage, but other I/O connects do have feature advantages and efficiencies that could be useful in bladed systems. There is also the distinct possibility that the partial physical layer convergence will continue and could yield physical unification in the 10Gb (XAUI) or 40Gb (XAUI DDR) generations.

This is quite an advancement rate. From pencil sketches in late 1998 to draft standards for high performance modules with full power thermal envelopes fully connected by meshes with multi-terabit class bandwidth built from mainstream commodity processors and I/O technologies in late 2002.



**On behalf of Hewlett-Packard, I'd like to thank the IEEE Computer Society and the Task Force on Cluster Computing for the opportunity to present to you today.**