# High Throughput Computing
## Linux Clusters and Grids

Carl G. Tengwall

IBM EMEA

tengwall@se.ibm.com

All statements regarding IBM's future direction and intent are subject to change or withdrawal without notice, and represent goals and objectives only.  Contact your IBM local Branch Office or IBM Authorized Reseller for the full text of a specific Statement of General Direction.

IBM may have patents or pending patent applications covering subject matter in this presentation. The furnishing of this presentation does not give you any license to these patents.

The information contained in this presentation has not been submitted to any formal IBM review and is distributed AS IS.

The following terms are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both: IBM, IBM logo, AIX, AIX/L, PowerPC, RS/6000, SP, Netfinity, pSeries, xSeries, Chipkill, ServeRAID.

Intel and Pentium are trademarks or registered trademarks of Intel Corporation in the United States, other countries, or both.

Myrinet is a trade name of Myricom, Inc.

LINUX is a registered trademark of Linus Torvalds

Other company, product, and service names may be trademarks or service marks of others.

# High Throughput Computing

## Agenda

➤ **Linux Clusters**

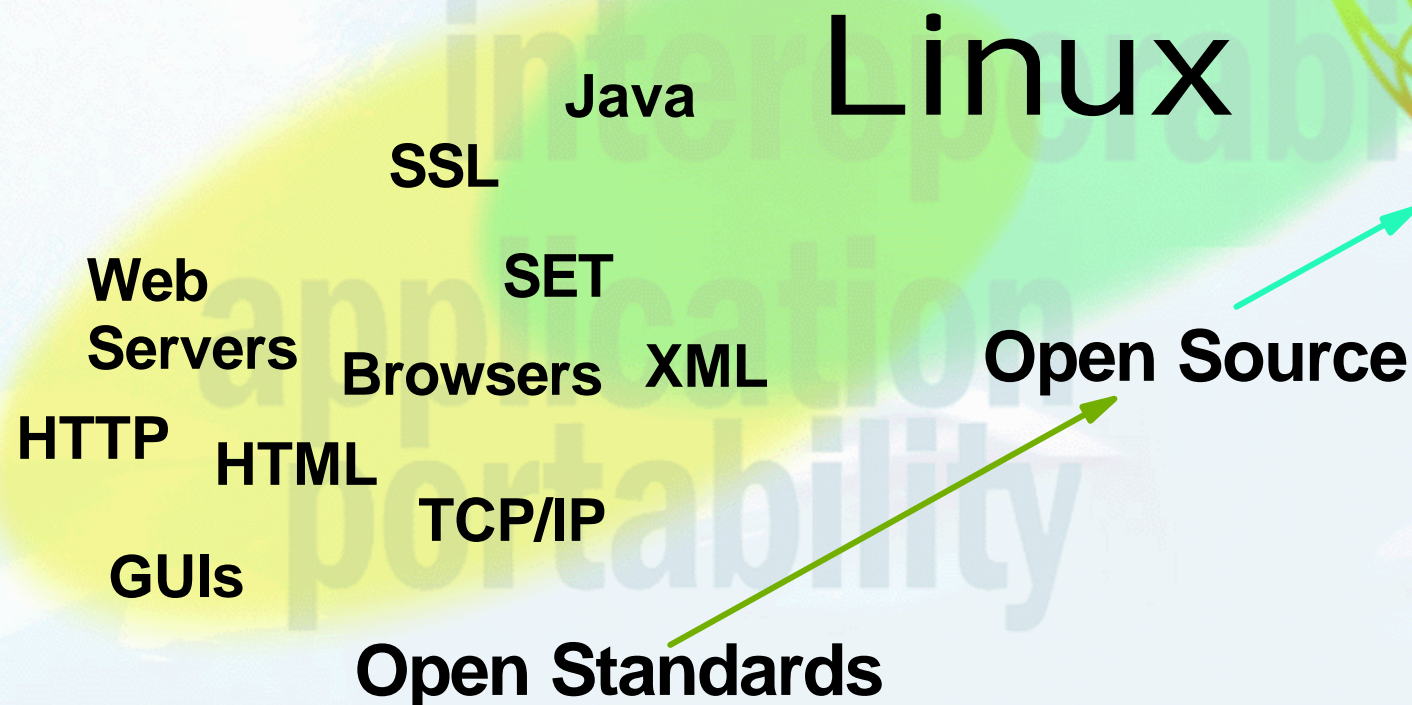➤ **Grids**

IBM

## Linux is important to IBM

- **Entrenched Internet technology**

- **Increasingly used in the HPC market**

- **Can become the volume Application Development and Deployment environment**

- **Potential to be a key technology for the next generation eBusiness**

**"Linux will do for applications, what the Internet did for networks"**

**Already #2 reference platform for application development**
- **Can become pervasive over time**

Next Generation eBusiness

## Linux

Java

SSL

SET

Web Servers

Browsers

XML

HTTP

HTML

TCP/IP

GUIs

**Open Source**

**Open Standards**

IBM

**Linux has real and perceived limitations today for pervasive, enterprise-wide use**

**IBM sees Linux as a strategic technology**

- **We are investing considerable resources and $$, and contributing key IBM technology to making it enterprise-ready**
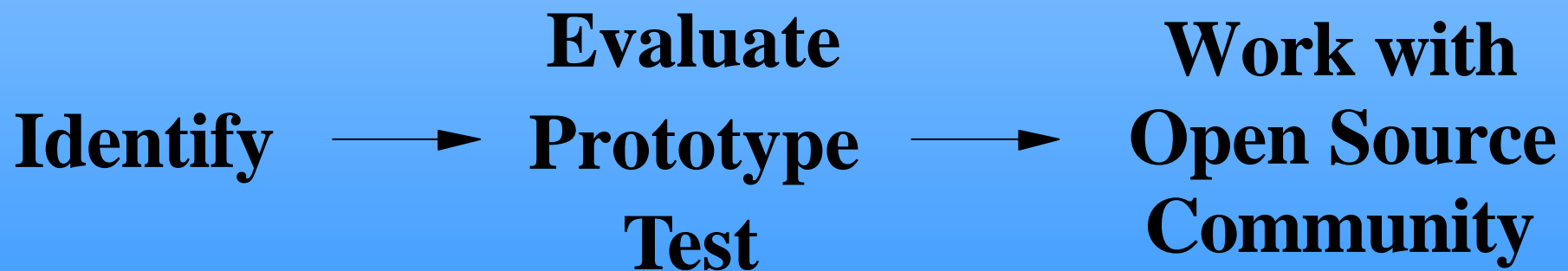- **We will work with the community to do that**

**Unix and other proprietary operating systems will continue to exist for foreseeable future**

- **Large investment in software, applications, data**
- **IBM continues to invest in AIX and pSeries systems**

# Linux Technical Strategy

- **Support Linux on all IBM platforms**

- **Strong affinity between IBM operating Systems and Linux**
  - **Example: AIX/L**

- **Work with the Linux community to infuse technology into the Linux kernel**

- **Deliver robust Linux Cluster solutions based on Open Source and IBM technologies**

- **Encourage adoption of Linux**

## Goal:   Accelerate maturation of Linux into Enterprise

- Distributed - worldwide organization of ~200 developers

- IBM's primary interface to open source Linux community

- Identify and work on enhancements for enterprise-class capability

**Identify** → **Evaluate Prototype Test** → **Work with Open Source Community**

IBM

**Linux Community - Core Components**

**Linux Technology Center**

Scalability - Resource and SMP
Journaled File System - JFS port to Linux
IA64 port - Project Trillian participation
Threads
Networking - Protocols, Device Drivers
Systems Management
Mathlib work - IA64 high-precision math functions
Linux Standards Base participation
Logical Volume Manager
File / Print Services
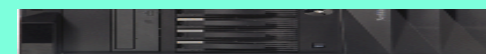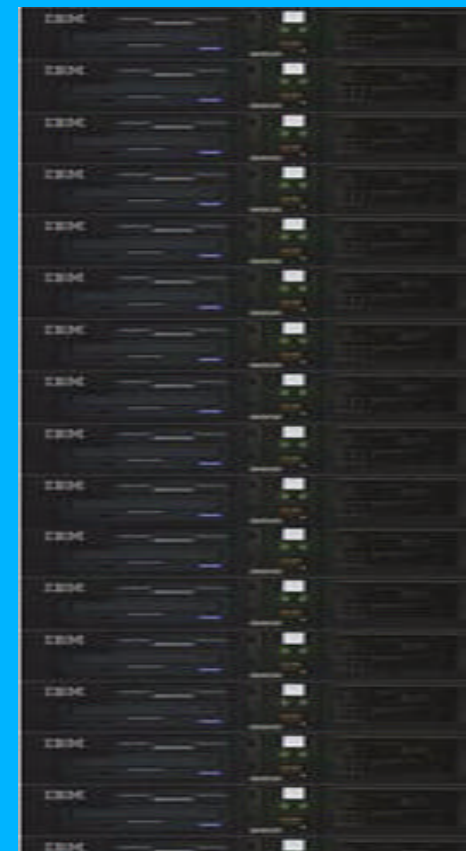SashXB - part of GNOME foundation technologies
GNOME foundation member

Distributions

**xSeries**  **pSeries**  **iSeries**  **zSeries**  **Appliances Storage**

IBM

# IBM Linux Cluster solution for S&TC

- **Prepackaged, prevalidated Linux Cluster**
- **1U and 3U 2-way IA-32 servers**
  - **PowerPC and IA-64 in '02**
- **Cluster and management networks, remote control**
- **Fully integrated, high availability storage solution**
- **Comprehensive Systems Management (CSM)**
- **Cluster File System (GPFS)**
- **S&TC optimized solution includes**
  - **High-performance Myrinet 2000 network**
  - **High-performance compilers (Fortran, C, C++, OpenMP)**
  - **Parallel Debugger (TotalView)**
  - **Job Management software (PBS)**
- **Optional enterprise service/support for IGS**

**2-way 1.26 GHz**

**4 GB RAM**

**146 GB disk**

# Rack-optimized IA-32 Systems

|  | xSeries 330 | xSeries 340 | xSeries 350 | xSeries 370 |
|---|---|---|---|---|
| Processor | 2-way SMP | 2-way SMP | 4-way SMP | 8-way SMP |
| Package | 1U | 3U | 4U | 8U |
| Max Memory | 4GB | 4GB | 16GB | 32GB |
| Internal HDD | 2 | 3 | 6 | 2 |
| PCI slots | 2 | 5 | 6 | 12 |

Integrated Service Processor
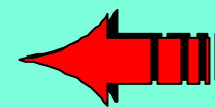
Netfinity Director

Software Rejuvenation

Processor Deallocation

Cable Chaining Technology

Predictive Failure Analysis
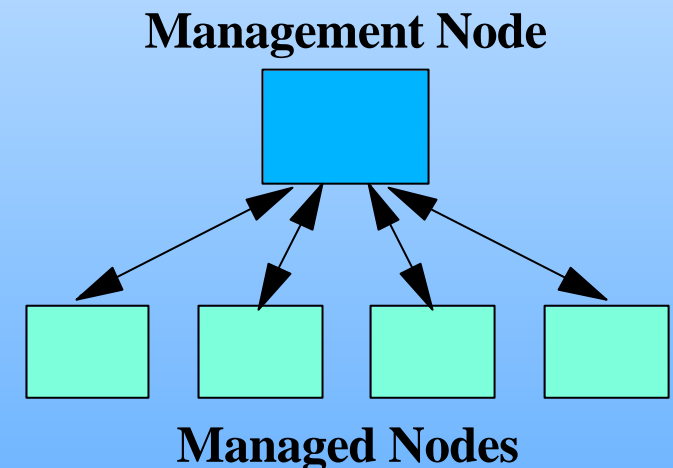
Hot Plug disk, adapters, fans, power

LightPath Diagnostics

ChipKill Memory

Varying capability in various models

IBM

# Cluster Systems Management for Linux

**CSM allows a cluster to be managed as a single entity from a single point of control**

- **Remote hardware control and monitoring**
  - **Power on/off/reset**
  - **Monitor environmental conditions**

- **Remote console function**
  - **Access to cluster servers prior to OS installation or when network access is unavailable**

- **Software installation**
  - **Cluster-wide parallel install**

- **Distributed Shell, Node Groups**
  - **Execution of arbitrary commands or scripts on all or some of the servers in the cluster**

**Management Node**

**Managed Nodes**

IBM

# Cluster Systems Management for Linux

**Allows a cluster to be managed as a single entity from a single point of control**

- **Configuration File Manager**
  - **Enables administrator to set up configuration files in a central place**
  - **An agent that pulls any changes down to each server in the cluster**

- **Distributed Management Server**
  - **Coordination for various management functions**
  - **Persistent repository of cluster configuration**
  - **Heartbeat function**
  - **Liveness state that can be assessed by other applications**

**Allows a cluster to be managed as a single entity from a single point of control**

- **Event Response Resource Manager**
  - **Mechanism for automatic response to specific events**
  - **Set of predefined events and actions that are commonly used in managing a cluster will be provided**
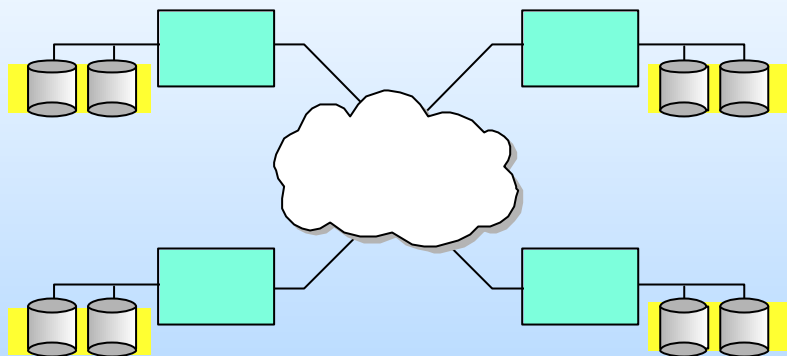
- **Probe manager**
  - **Set of probes to check consistency of cluster configuration information and diagnose configuration errors**

**Much of CSM is based on mature SP technology**
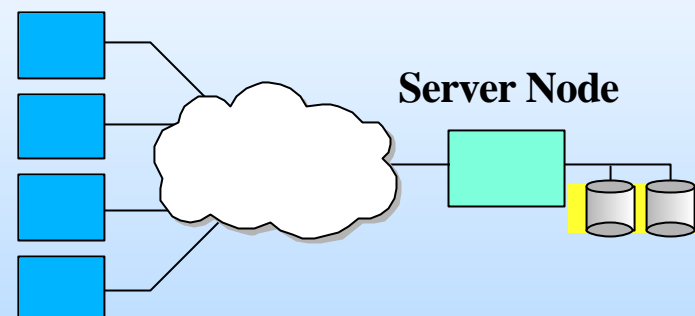**Used in over 10,000 SP systems today**
**Gone through multiple releases over past 9 years**
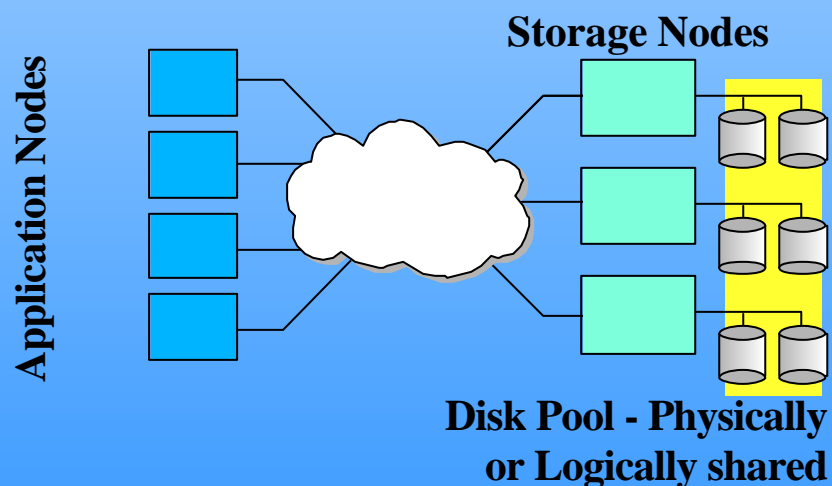
# Scalable Cluster File System

- **Native File System**
  - **No file sharing - application can only access files on its own node**
  - **Applications must do their own data partitioning or replication**

- **DCE Distributed File System**
  - **Application nodes share files on server node**
  - **Coarse-grained (file or segment level) parallelism**
  - **Server node is performance and capacity bottleneck**

**Application Nodes (Clients)**

**Server Node**

**Storage Nodes**

**Application Nodes**
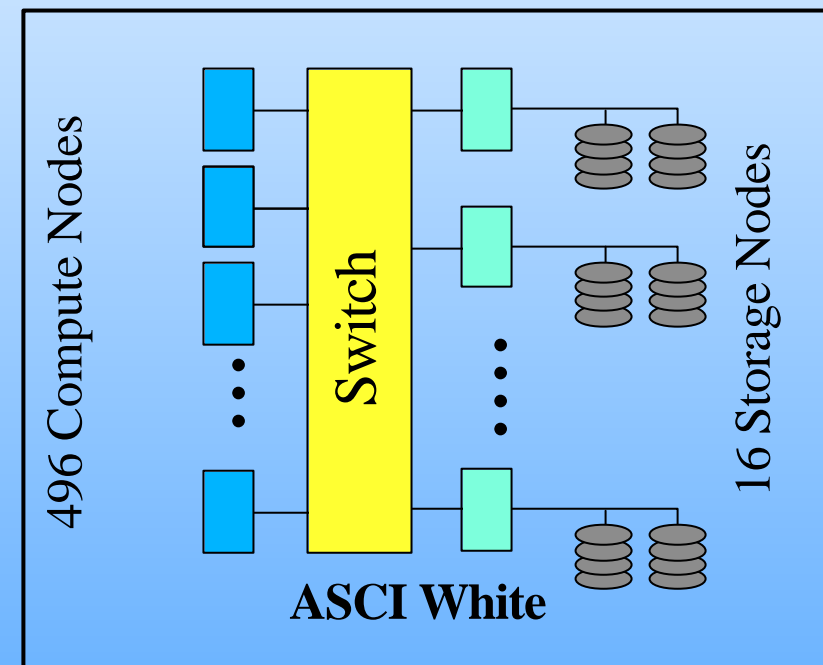
**Disk Pool - Physically or Logically shared**

- **GPFS Parallel File System**
  - **Striped across multiple disks on multiple storage nodes**
  - **Independent GPFS instances run on each application node**
  - **Storage nodes used as "block servers"**
    - **all instances can access all disks**

IBM

- **Posix standards-compliant**

- **Uniform access via (logically or physically) shared disks**

- **High capacity - tens of TB per file system, ~TB per file**

- **High throughput**
  - **Wide striping and large data blocks**
  - **Client caching via distributed locking**
  - **Parallel access via fine-grained (byte range) locking**
  - **High sequential throughput via aggressive prefetch**

- **Reliability and fault-tolerance - node and disk failures**
  - **Journaling, data replication, RAID support**
  - **High-availability infrastructure**

- **Export via NFS and DFS**

- **GPFS is in it's fifth release and is used in widely used by our RS/6000 SP customers (AIX)**

- **ASCI White**
  - **512 nodes, 8K CPUs**
  - **150+ TB**
  - **12 GB/s to/from single file (or multiple files)**

- **Can be used as scalable NFS or DFS**

- **Now available on IA32/Linux**

**496 Compute Nodes** | **Switch** | **16 Storage Nodes**

**ASCI White**

# IBM Global Services

## Education & Training

- **Classroom or via web**
- **Available in 20 countries, multiple languages**
- **How-to (Redbooks) for Linux**

## Service & Support

- **24 X 7 enterprise level support**
- **All major distributions**

## Professional Services

- **Comprehensive enterprise services for Linux**
- **Infrastructure consulting and planning**
- **Installation**
- **Configuration**
- **Application enablement**

- **Established April 2001**

- **Mission:  To provide education and advanced technical training for the deployment and use of Linux based computing clusters for the HPC community, worldwide**

- **Three key players:  NCSA, AHPCC, IBM ACTC**

- **Intensive hands-on workshops, System Admin training, application tuning**

- **First workshop Oct 1-5 2001 at NCSA**

- **Worldwide workshops in 2002**

- **http://www.linuxclustersinstitute.org**

## IA-64 and PowerPC Linux

## Scalability optimizations

## Functional Enhancements
- Install
- Automated operations
- Security
- Usability

# Linux Cluster Successes

## U of New Mexico

- 256 x330s
- 80th on the 12/00 Top500 Supercomputers list

## Maui High Performance Computing Center

- 288 x330s
- One of the larger SP sites

## NCSA

- Support next generation Grid
- xSeries servers: 512 x330s and 100+ IA-64 nodes (1 TF each)
- IBM SW to support scaling, management and application in a tera-scale Linux cluster environment

## Royal Dutch Shell

- Tera-scale seismic processing
- 1024 x330s (1+Tflop)
- IBM Global Services to design, build, and implement

## MDS Proteomics

- Two 100-node x330 clusters
- 80th on the Top500 Supercomputers list

## weather.com

- One of "top 25" web sites
- xSeries servers with Linux, Websphere Commerce Suite, IBM Global Services design approach
- Cost, availability, scalability requirements

## Agenda

➤ **Linux Clusters**

➤ **Grids**

IBM
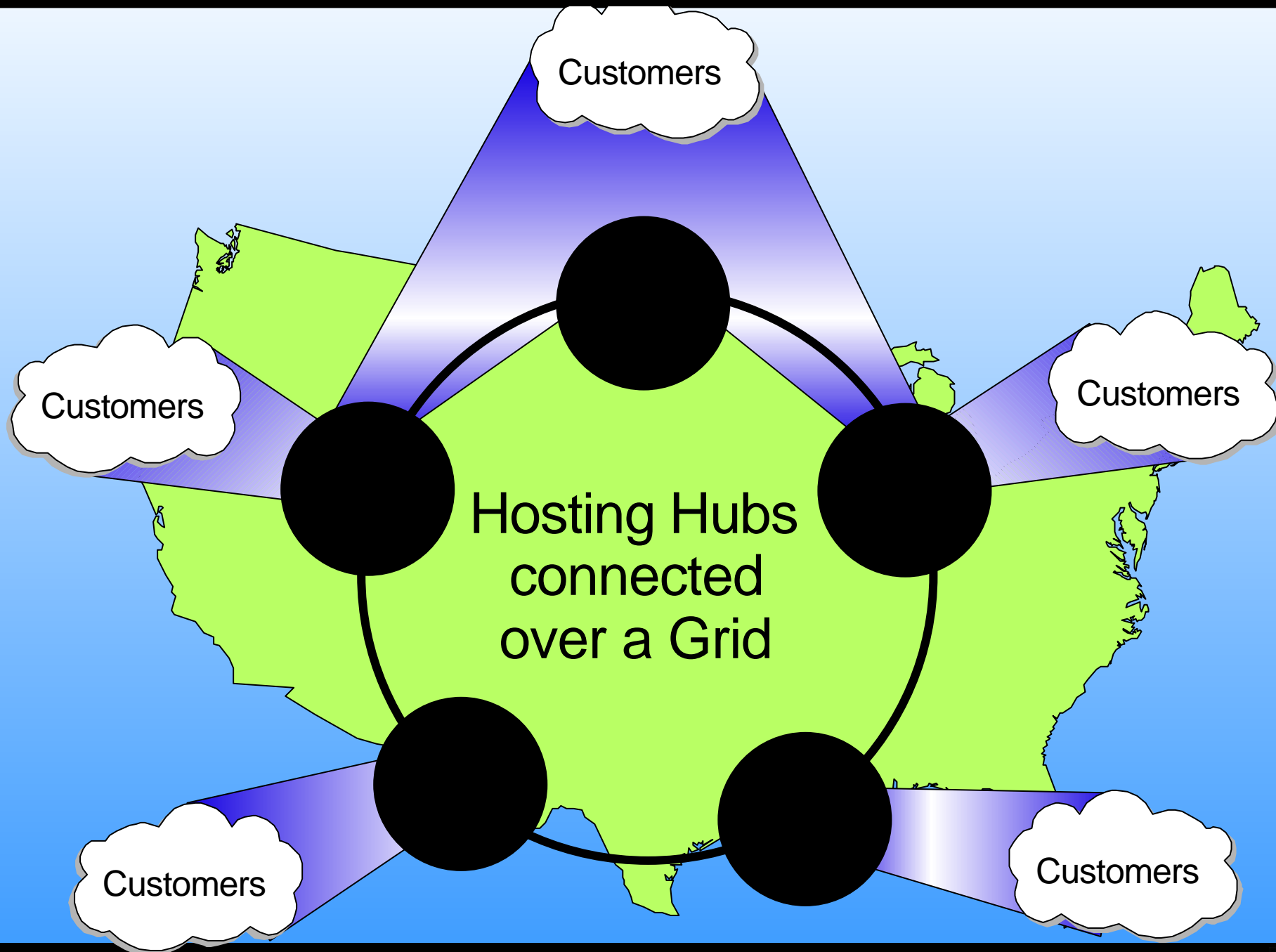
- **We believe Grids will emerge much as eSourcing and eServices will**

- **We recognize Grid Computing as a key strategic area**

- **Irving Wladawsky-Berger leads the IBM Grid Computing Initiative**
  - **Forming a cross-unit design council**
  - **Align with eSourcing strategy**
  - **Engage with Grid development community**
  - **Encourage joint University research in appropriate areas**

- **Delivery of standardized processes, applications, and infrastructure over the network as a service on a pay-as-you-go basis**
  - **Business Functions: CRM, eCommerce, Supply Chain, ...**
  - **IT Functions: Security, Web Hosting, Storage Services, Systems Management, ...**

- **$4B to add 50 hosting centers worldwide to serve as eSourcing hubs**

- **First step in the Utility model**

**Connect hubs into a Grid over time**

# eSourcing

Customers

Customers

Customers

## Hosting Hubs connected over a Grid

Customers

Customers

IBM

"Although other companies have expressed interest [in Grids], Foster and Hey said IBM has shown the most significant support so far. 'IBM is distinguished by farsightedness and enthusiasm,' Hey said.  'This stuff, to be significant in the long term, has to move into the commercial space, and IBM has stepped up,' Foster said."

**New York Journal News, Aug 2.**

**Work with the community**

**IBM technology where relevant**

**Grid-enable IBM products**

**Promote use in a wider segment**

**Very similar to the IBM Linux Initiative**

# Challenge

- ► **Build the world's most powerful computing GRID for consortium of major research facilities**

- ► **Create the "TeraGRID" infrastructure integrating fastest supercomputers, high-resolution environments, toolkits, and data storage facilities**

# Solution

- ► **Interconnected series of Linux clusters able to process 13.6 trillion calculations per second using open protocols**

- ► **Four DTF sites enabling thousands of scientists to share resources over the world's fastest research network**

**National Center for Supercomputing Applications (NSCA)**

**San Diego Supercomputing Center (SDSC)**

**Argonne National Laboratory**

**California Institute of Technology**

# Concluding Remarks

- **Linux is a key component of IBM's strategy. We are following a multi-pronged approach to accelerate adoption of Linux into the Enterprise**
  - Support Linux on all IBM platforms
  - Build strong affinity between IBM operating Systems and Linux
  - Work with the Linux community to infuse technology into the Linux kernel
  - Facilities to assist migration to Linux
  - Deliver robust Linux-based solutions using Open Source and IBM technologies

- **We are focusing on developing and deploying technology that will make the configuration, management, and efficient use of Linux systems and Linux Clusters easier in the Enterprise**

- **Grid computing intersect several key IBM initiatives and strategies. We will work with the community to define and deploy a robust infrastructure and accelerate its adoption across a wider segment**

IBM

Server Group

**ibm.com/linux**

**ibm.com/developerworks**

**oss.software.ibm.com/ developer/opensource/linux/**

# Thank you!

## Questions?

Linux at IBM