



**LINUX
NETWORKX™**

Linux NetworkX HPC Strategy and Roadmap

Eric Pitcher

October 2006

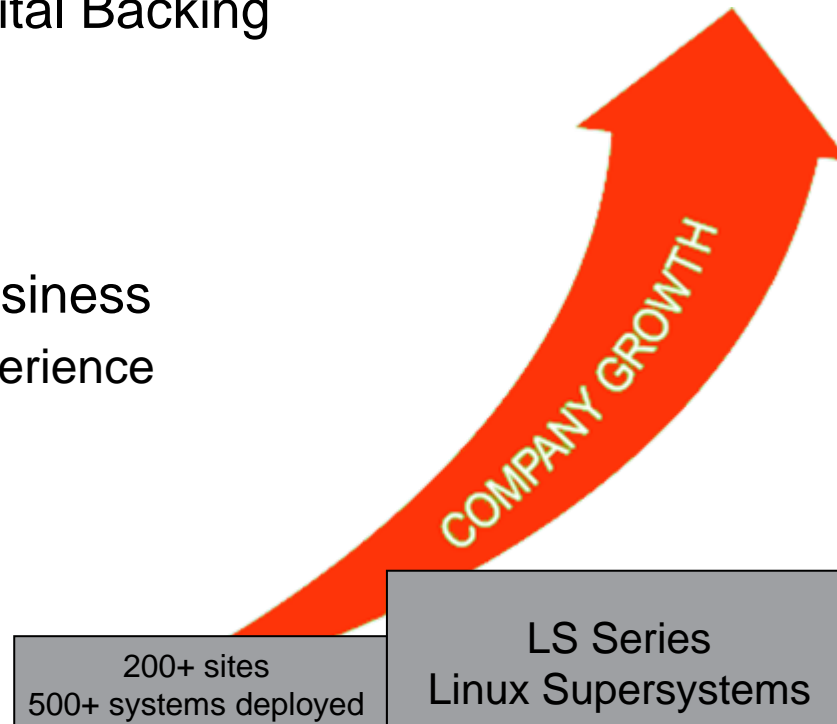
The Linux Supercomputing Company™

Agenda

- Business Update
- Technology Trends
- Linux Network Drivers
- Hardware Roadmap
- Software Highlights

Linux Network Overview

- Founded in 1989, HQ in Salt Lake City
- Operations in Americas, EMEA, Asia Pacific
- Privately Held – Strong Venture Capital Backing
 - Oak Investment Partners
 - Tudor Ventures
 - Lehman Brothers
- Linux Supercomputing is our only business
 - Nearly a decade of Linux cluster experience



Selected Linux Network Large Systems



Site: US Army Research Laboratory
Computer Name: John Von Neumann
Theoretical Peak: 13.926 TFlops
Best Top500 Ranking: 13
System Model: Evolocity II
Processors: 2048 Xeon 3.4 GHz



Site: Los Alamos National Laboratory
Computer Name: Lightning
Theoretical Peak: 11.26 TFlops
Best Top500 Ranking: 6
System Model: Evolocity
Processors: 2816 Opteron 2 GHz



Site: Lawrence Livermore National Lab.
Computer Name: MCR
Theoretical Peak: 11.2 TFlops
Best Top500 Ranking: 3
System Model: Evolocity II
Processors: 2304 Xeon 2.4 GHz



Site: Los Alamos National Laboratory
Computer Name: Lightning Bolt
Theoretical Peak: 16.645 TFlops
Best Top500 Ranking: N/A
System Model: Evolocity II
Processors: 1536 Opteron 2.4 GHz
+ 245 Dual Core 1.8 GHz



Site: Los Alamos National Laboratory
Computer Name: Pink
Theoretical Peak: 10.0 TFlops
Best Top500 Ranking: N/A
System Model: Evolocity
Processors: 2050 Xeon 2.4 GHz



Site: NERSC
Computer Name: Jacquard
Theoretical Peak: 3.1 TFlops
Best Top500 Ranking: N/A
System Model: Evolocity II
Processors: 722 Opteron 2.2 GHz

Selected Recent Orders

- **Fleet Numerical Meteorology and Oceanography Center**
 - 229 Nodes & 4 Accelerator Nodes – mid-Nov
- **DoD TI-06**
 1. Classified System – 842 Intel Dempsey Nodes
 2. Non-classified System – 1024 Intel Woodcrest nodes
 3. Dugway – 64 node Dempsey
 4. ARL/Visualization – 64 node
 5. TDS – 16 nodes Woodcrest/Dempsey
- **NASA/Goddard**
 1. Baseline System – 128 Dempsey nodes
 2. Visualization System – 16 node
 3. SU01 – 256 Woodcrest Nodes
 4. SU02 – 256 Woodcrest Nodes

What Drives Us

**Create the best Linux Supercomputing systems,
software and services**

to help our customer's solve important problems.

Technology Trends

Processor Highlights

CPU

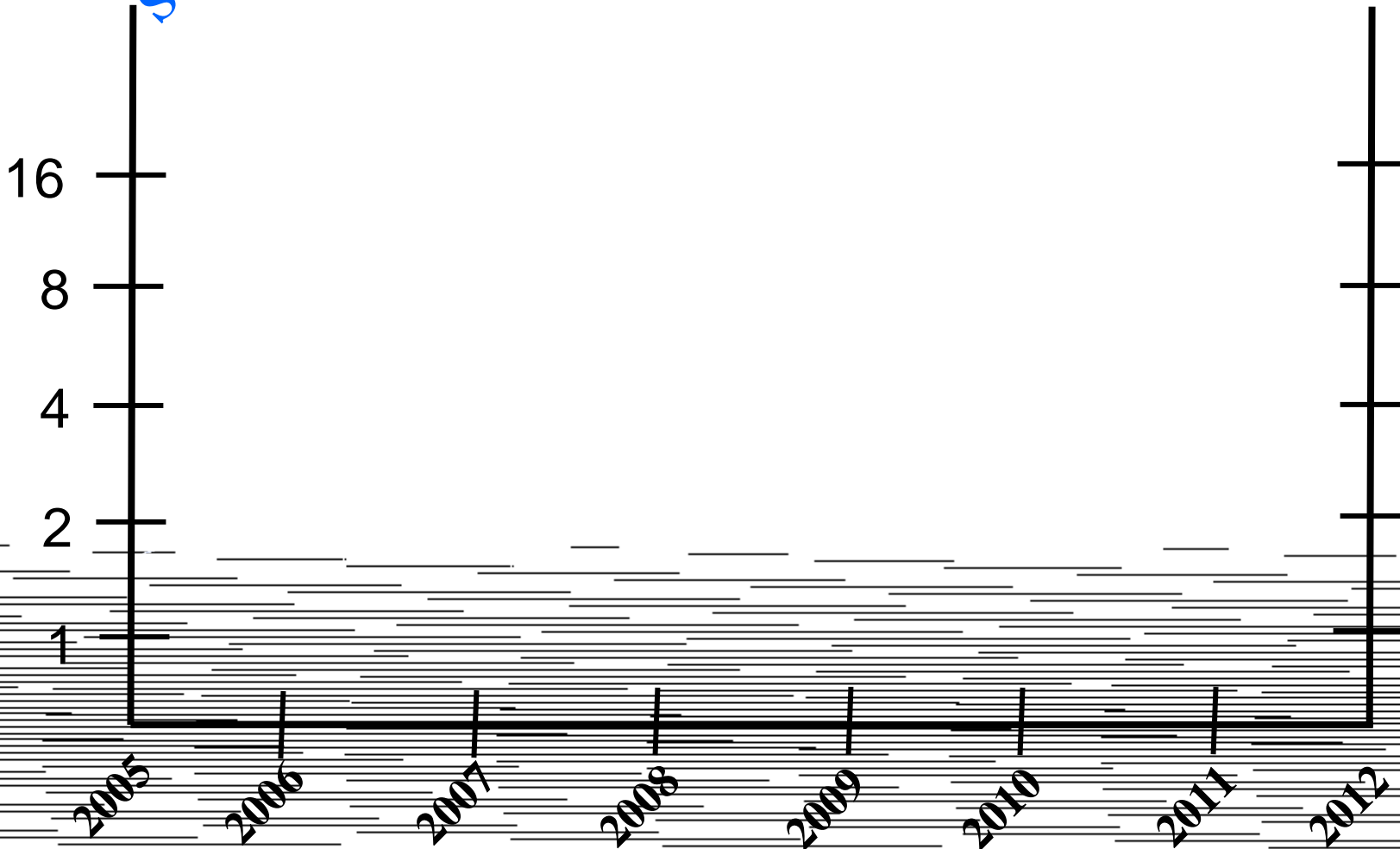
- On-die memory controllers seem to be a likely direction (improved latency)
- Trend towards point-to-point FSB (improved on-node scalability)
- Both major architectures going towards 4 FLOPS/clock
- Focus on power efficiency through 2012

Multi-paradigm Computing

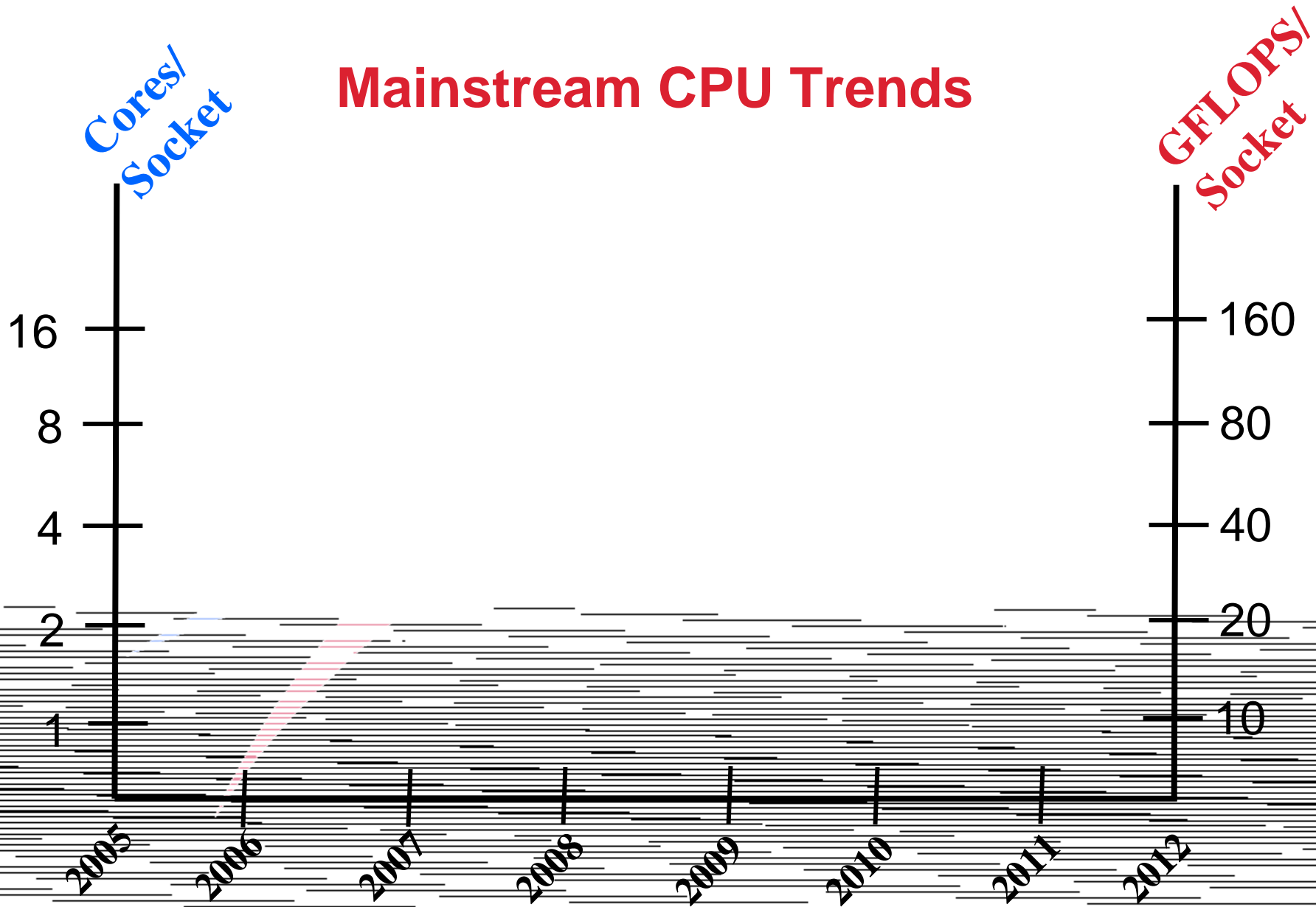
- Increase in the use of “co-processors” for multi-paradigm computing
- Multi-paradigm computing will provide highly capable offload engines -- initially difficult to program

Mainstream CPU Trends

Cores/
Socket



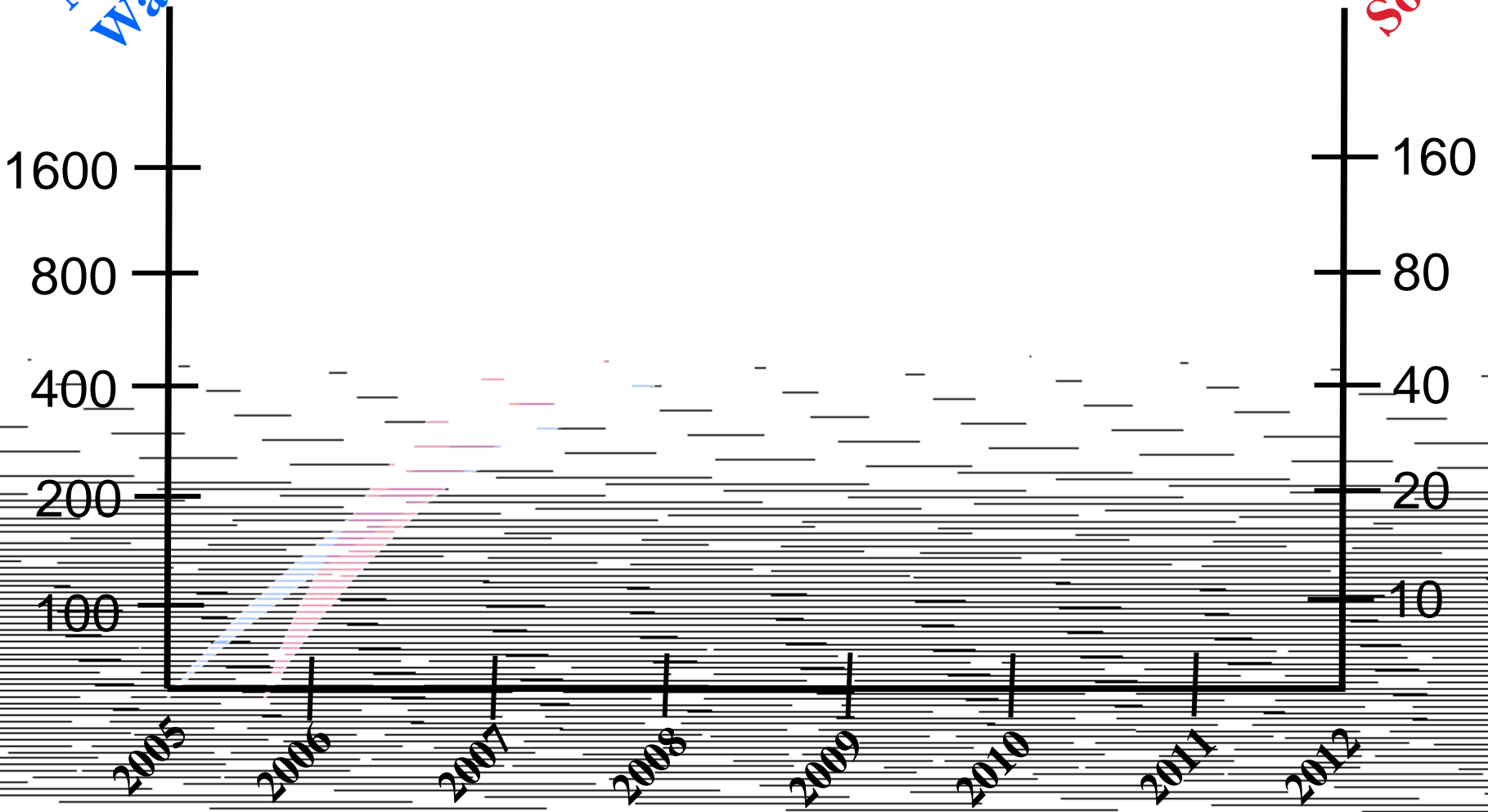
Mainstream CPU Trends



CPU Power Trends

MFLOPS/
Watt

GFLOPS/
Socket



CPU Summary

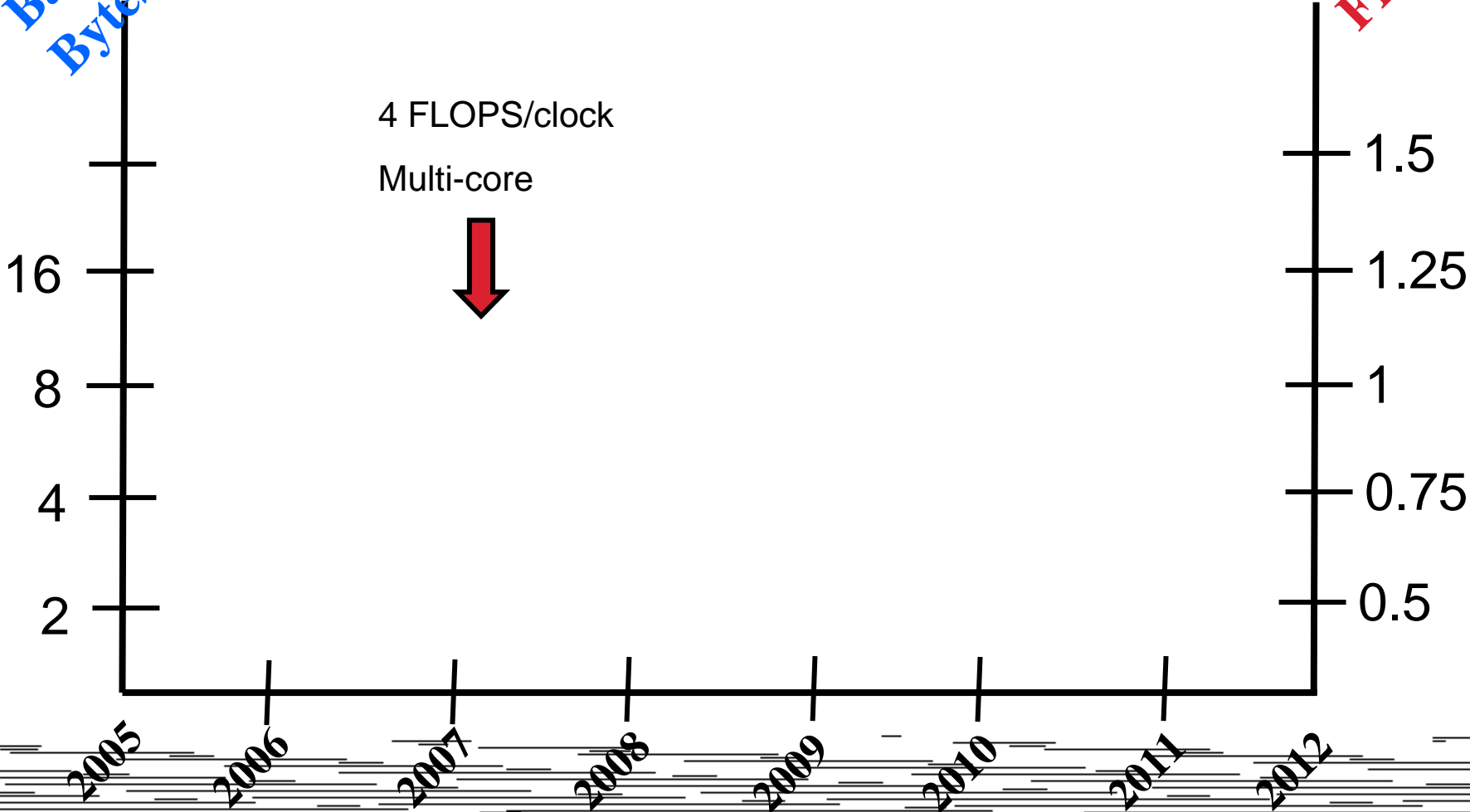
- Multi-core CPUs will greatly exceed single-core CPUs in computing power
- The cost-effectiveness of multi-core CPUs may make it attractive to apply fewer cores to an application to get more bandwidth.
- Renaissance in CPU and co-processor architectures expected over the next 5 years. Expect to see a lot of interesting ideas. May be hard to pick winners.

Memory Bandwidth

Per-core
Bandwidth
Bytes/clock

Bytes/
FLOP

4 FLOPS/clock
Multi-core



Memory Summary

- Bytes/FLOP stays mostly constant through 2010. So performance per node should significantly increase.
- Remedy for memory-intensive apps: turn off some cores leaving more bandwidth for remaining cores.
- Latency is currently inferior for FB-DIMMs, but capacity is better.

Interconnect Highlights

- Infiniband will offer leading-edge performance through 2009
 - Bandwidth will increase rapidly through 100 Gbps
 - Latencies will likely flatten out @ ~500-800 ns by 2011
 - Significant improvements in optics in the near future
- 10Gb Ethernet will become more cost effective

Infiniband in 2007

- Half-round trip ping-pong latency < 1.5 μ s
- Messaging rate > 10 M/sec
- Inexpensive optical cabling for SDR and DDR up to 100 m
- More mature stack natively supported in SLES 10
- More goodies TBA later this year

Linux Networx Drivers

Key Engineering Themes

- Use “Off the shelf” to optimize price/performance
- Systems
 - “Standard,” supported systems
 - Integrated, validated, tested software stack
 - Innovate “on top” of standard hardware and software
- Environment/Pricing
 - TCO increasingly important (eg, quick system deployment)
 - Power and cooling requirements becoming more important
- Focus on production supercomputing
 - Necessitates full-featured software stack

LNXI Systems Strategy

2007 Cluster System Mgmt Software

- Performance & Utilization
- System Usability & Management
- RAS
- Operable on other vendor systems

2006 Application-tuned Supersystems

- Performance-tuned for specific applications
- Production at Power-up

2006 Performance Software Platform

- Integrated/validated software stack

2005/6 Standard Hardware Platforms

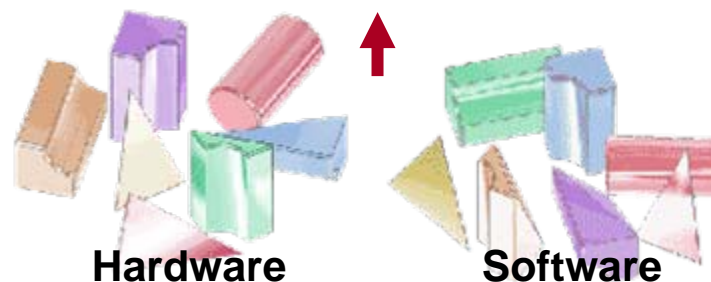
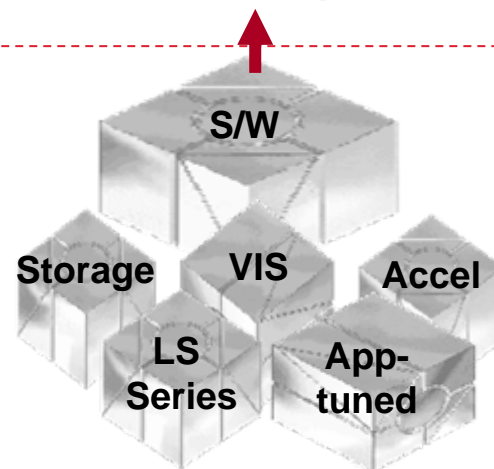
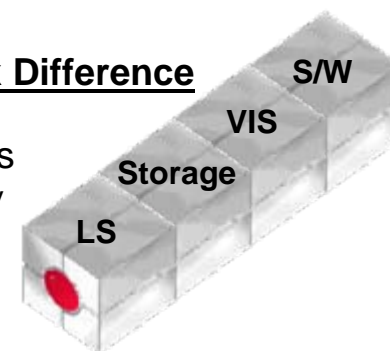
- Integrated systems delivery

1st Generation Clusters

- Pre 2006 - LNXI leadership with stable “systems”
- SC expertise and credibility
- Custom approach less scalable

The Linux Network Difference

- Application Expertise
- Integrated Systems
- Management/usability
- RAS
- Performance
- Utilization

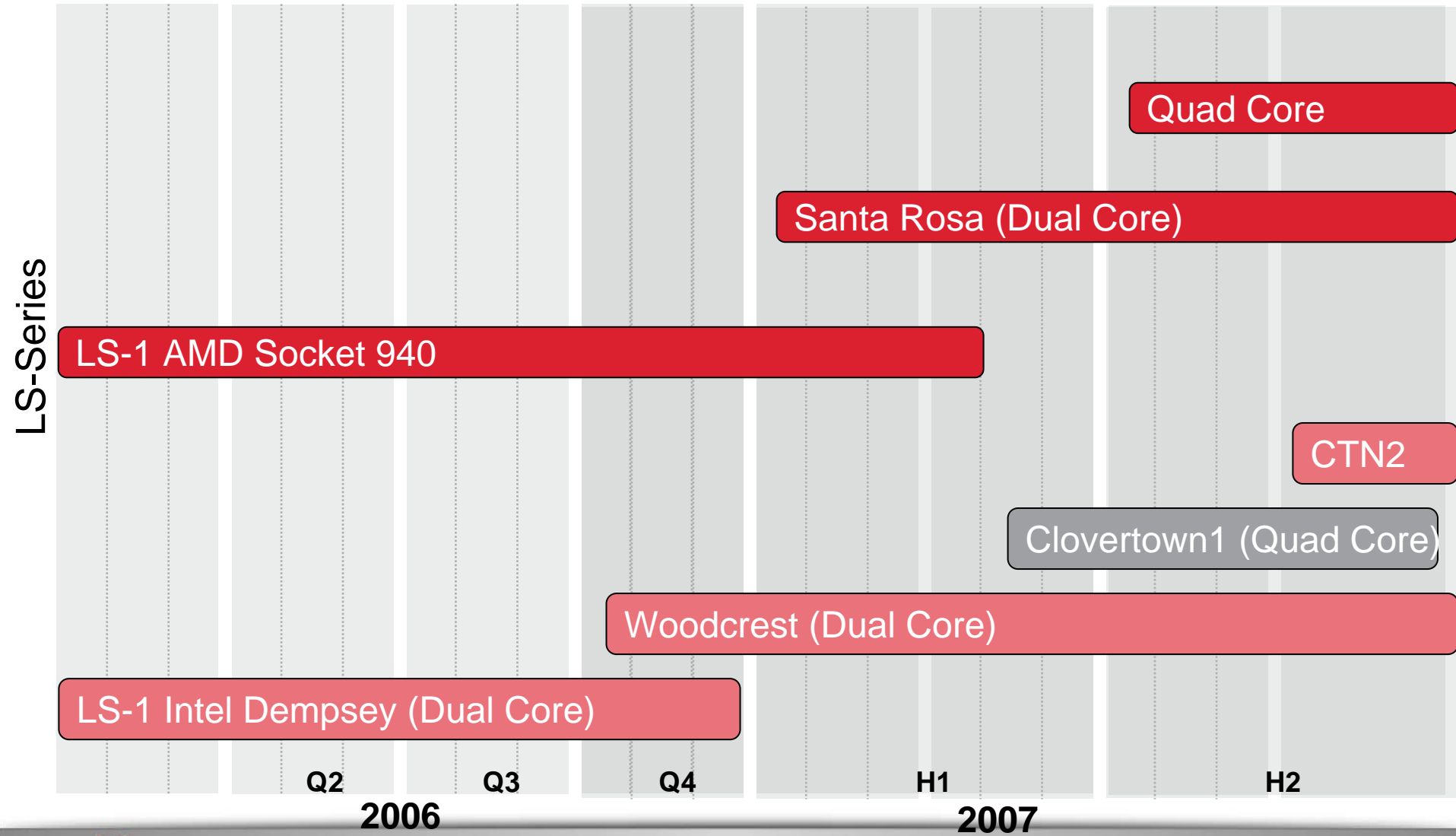


Hardware Roadmap

LS Series

LS-1 AMD

LS-1 Intel



Value Performance

**Total Cost of Ownership
No In-house S/W Development
Limited Linux Experience**

Small File
I/O

NFS

Premium Performance

**Multiple applications/balanced system
Online capacity expansion
Production system environment**

Diverse Streaming
and Small File
Applications

CFS Lustre
Linux Network
GPFS

Ultimate Performance

**Extreme I/O performance
Savvy In-house technologists
Time to Results is top priority**

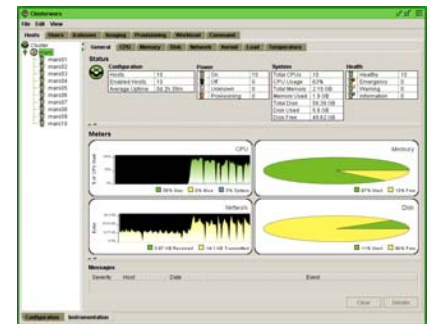
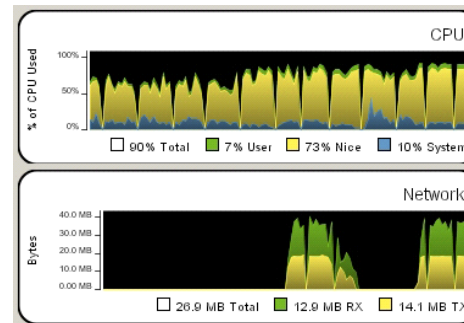
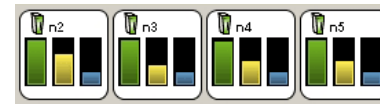
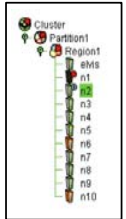
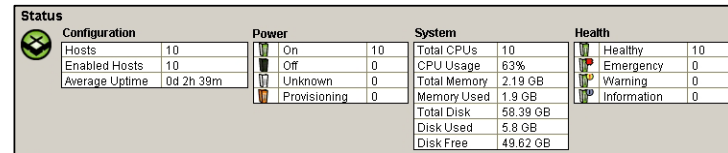
Streaming
Bandwidth

Software Highlights

Key Linux Networx System Features

Clusterworx

- Maximizes system performance with comprehensive system monitoring.
- Increases uptime with automated system management.
- Quickly updates the entire system with fast multicast provisioning.
- Implements risk-free changes to OS, applications, and kernels through version-controlled image management
- Allows users to assess the state of the entire system and each node at a glance.



System Dashboard

Clusterworx - cwxhost

File Edit View

Hosts Users Iceboxes Imaging Provisioning Runner

Cluster Host Nodes-Highm Nodes-Lowm

cuda1 cuda2 cuda3 cuda4 cuda5 cuda6 cuda7 cuda8 cuda9 cuda10 cuda11 cuda12 cuda13 cuda14 cuda15 cuda16

cuda17 cuda18 cuda19 cuda20 cuda21 cuda22 cuda23 cuda24 cuda25 cuda26 cuda27 cuda28 cuda29 cuda30 cuda31 cuda32

General CPU Memory Disk Network Kernel Load Temperature

Configuration

Hosts	16
Enabled Hosts	16
Average Upt...	4d 4h 55m

Power

On	8
Off	8
Unknown	0
Provisioning	0

System

Total CPUs	14
CPU Usage	100%
Total Memory	62.4 GB
Memory Used	2.77 GB
Total Disk	846 GB
Disk Used	42.7 GB
Disk Free	760 GB

Health

Healthy	0
Emergency	0
Warning	0
Information	16

CPU

Memory

Network

Disk

Messages

Severity	Host	Date	Event
Warning	cuda9	06/25/2006 04:22:22 AM	Clusterworx has detected that LinuxBIOS is running ...
Warning	cuda10	06/25/2006 04:02:06 AM	Five minute load average limit 1.1 exceeded on hos ...
Warning	cuda10	06/24/2006 04:08:23 PM	Clusterworx has detected that LinuxBIOS is running ...
Warning	cuda11	06/24/2006 04:03:10 AM	Five minute load average limit 1.1 exceeded on hos ...

Clear Details

Instrumentation Configuration

Key Linux Networx System Features (Cont'd)

- Full hardware validation
- Complete System Integration and Testing before Shipping
- Project Management
- Professional Services
 - Site Planning and Environmental Design
 - Application Parallelization Consulting
 - Data Storage and Management Assessment
 - Data Storage Design and Implementation
 - Capacity Planning
- Linux Networx Training
 - Provides users with information, tools and practical experience to successfully manage LNXI Supercomputing Systems.

WRF Model on LS-1 System Scaling: 64-512 Cores

