# Reflections on Building a High-performance Computing Cluster Using FreeBSD

NYCBUG
March 2008

Brooks Davis
The Aerospace Corporation
El Segundo, California, USA
http://people.freebsd.org/~brooks/pubs/meetbsd2007/

**THE AEROSPACE CORPORATION**

# Outline

- Fellowship Overview
- Evaluation of Design Issues
- Lessons Learned
- Thoughts on Future Clusters
- Conclusions

# A Brief History of Fellowship

- Started in 2001
- Primarily motivated by the needs of the GPS program office
- Intended to be a corporate resource
  - Diverse set of users
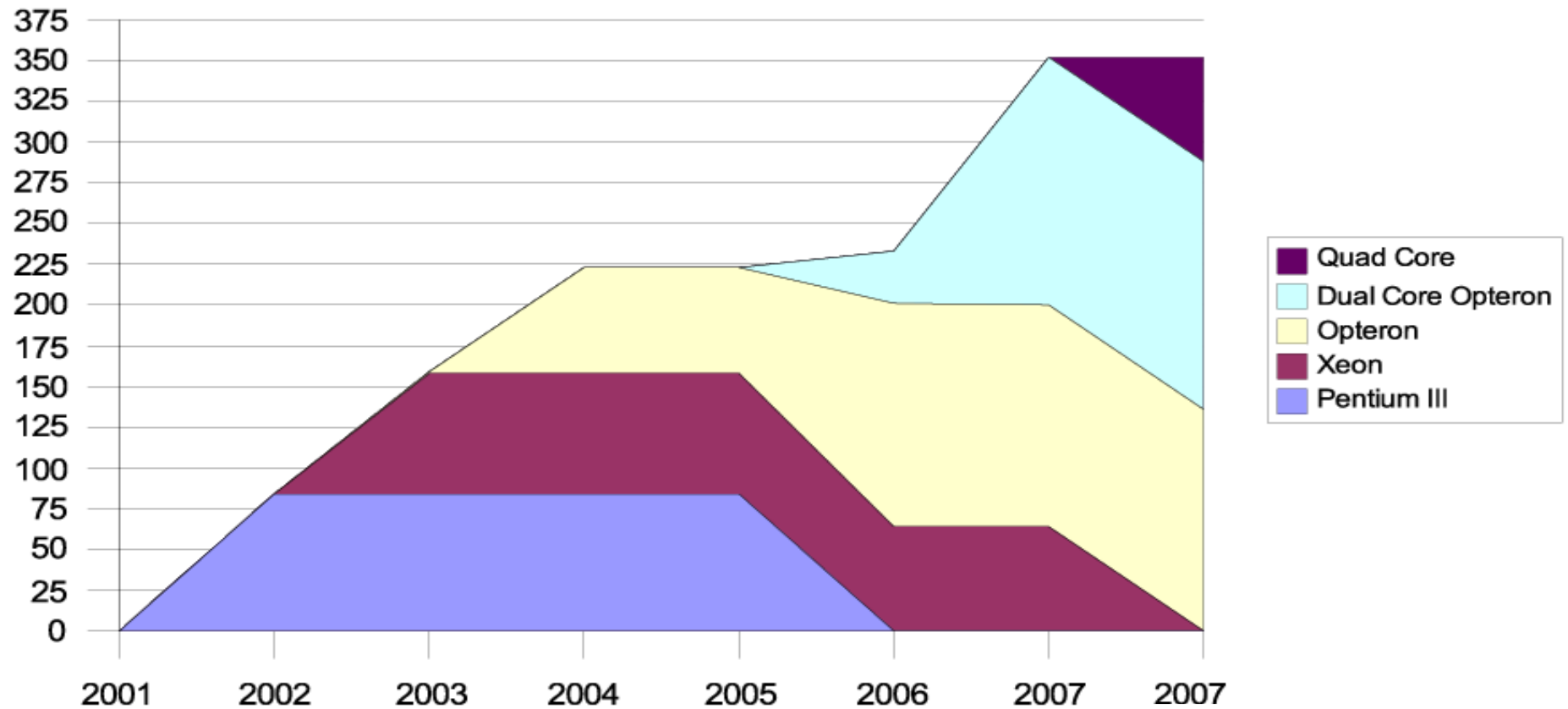
**THE AEROSPACE CORPORATION**

# Fellowship Hardware

- 352 dual-processor nodes
  - 64 Quad-core Woodcrest Xeons
  - 288 Opterons (152 dual-core)
  - 1-4GB RAM, 80-250GB disk
- 6 core systems
  - fellowship – shell server
  - fellowship64 – amd64 shell server
  - arwen – node netboot server, scheduler, NIS, DNS, Mathematica License Server
  - elrond – `/scratch`
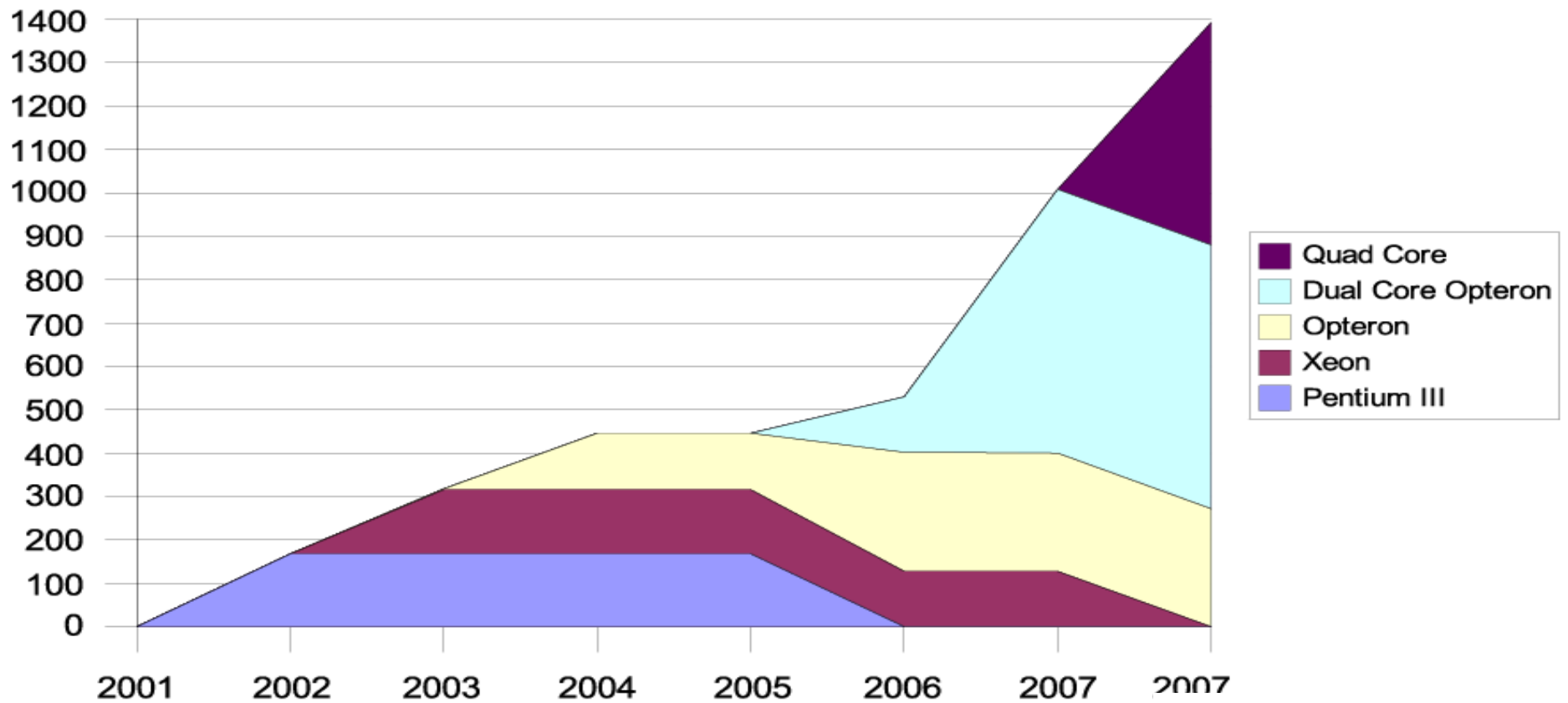  - moria – NetApp FAS250 `/home`, `/usr/aero`
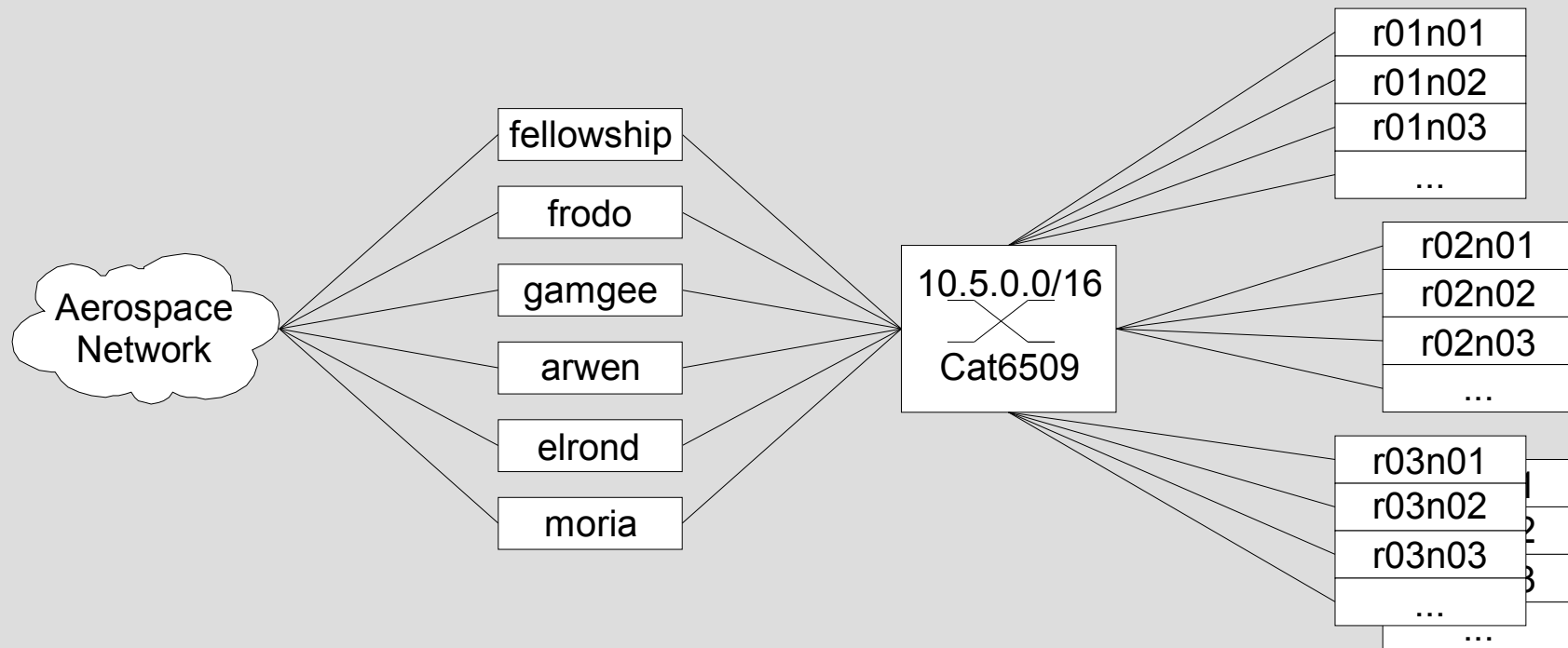
# Fellowship Circa February, 2007

# Fellowship Composition by processor count

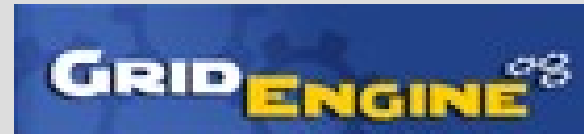# Fellowship Composition by core count

# Fellowship Network

# Fellowship Software

- FreeBSD 6.2-RELEASE-p8
- Sun Grid Engine (SGE) 6.0u11
- Message Passing Interface (MPI)
  - MPICH
  - MPICH2
  - OpenMPI
- PVM
- Mathematica
- Matlab

# Design Issues

- The Issue
- What We Did
- How it worked

# Operating System

- FreeBSD 4.x initially
  - Many years of service
  - Userland threading forced some application design choices
  - Extensive modification to boot scripts required
- FreeBSD 6.2-RELEASE now
  - Working well
- FreeBSD 7.x soon
  - DTrace in 7.1!

THE AEROSPACE CORPORATION

# Node Architecture

- Dual processor
  - Intel Pentium III initially
    - Worked well
    - Some system problems (unrelated to CPU)
  - Intel Xeon next
  - AMD Opteron
    - single-core -> dual-core
  - Intel Xeon "Woodcrest" today
- Netboot
- Disk for scratch
  - Destroy contents on crash
- Gigabit Ethernet
- Custom 1U rack mount chassis

THE AEROSPACE CORPORATION

# Form Factor

- 14in, 1U rack-mount chassis for Pentium III
- Xeons in standard 1U chassis
- Custom, front-port 1U chassis for Opterons

# Node Configuration Management

- Netboot with NFS roots
- 4.x: upgrade images in chroot
  - Eases maintaining `/etc` customizations
  - Makes it easy to forget `/etc` customization
- 6.x: NanoBSD derived scripts to build from scratch
  - Records modifications to `/etc`
  - Requires build documentation of unscriptable modification

# Physical System Management

- Initial setup
  - KVM switch for core systems
  - Serial consoles for nodes
    - BIOS access disabled due to hangs at boot
    - Rarely used
  - Cart with KVM for node maintenance
- Now
  - VNC based remote KVM switch for core
  - Rack mount KVM in new node racks for maintenance

# Lessons Learned

- Nothing shocking
- Types
  - Technical
  - Non-technical (user related)

THE AEROSPACE CORPORATION

# Uncommon Events May Become Common

- Examples:
  - Relatively rare (~1/30) BIOS hangs with serial redirection
    - Major issue with every downtime
  - PXE boot failures
  - Disk failures
  - Power supply failures
- Failures may be systemic
  - Disks were "DeathStars"
  - Power supplies clearly defective

# Neatness Counts

- Poor cable management, etc makes replacement more difficult
  - Poor cable planning
  - Sloppy installation
- Good argument for blades
  - Particularly for environments with minimally trained hands

# All the World is Not a Linux Box (but some vendors think it is)

- Some applications required porting
  - SGE
  - Ganglia
  - OpenMPI
  - GridMPI
  - Globus Toolkit
- Some applications are Linux only
  - Mathematica
  - Matlab
  - Total View
- Sometimes hard to tell what causes a failure
  - LAM-MPI

# System Automation is Critical

- Doing anything by hand on 352 nodes is impractical
- Without automated monitoring, users find problems the hard way
  - applications die, etc.
- UNIX tools to the rescue
  - xargs is your friend

THE AEROSPACE CORPORATION

# onallnodes script

```sh
#!/bin/sh
FPING=/usr/local/sbin/fping
NODELIST=/usr/aero/etc/nodes-all

${FPING} -a < ${NODELIST} | \
    xargs -n1 -J host ssh -l root host $*
```

# User Related Lessons

- In our environment users have a job to do and are experts, but:
  - It is generally not computer science
  - They often have a limited grasp of software engineering tools and principles
    - "You can write FORTRAN in any language"
- Users do not want to learn new techniques
  - Forcing them to use the scheduler after several years of it being optional was painful for everyone
- Extremely uneven user demands
  - Top user ~50% of cycles , top 5 >95%

# Thoughts on a Future Cluster

- Fully diskless
  – Disks were a top cause of failure
- Higher bandwidth, lower latency network
- Consider blades
- Reconsider OS choice
- Run by ordinary system administrators
  – Not researchers

# Conclusions

- Fellowship is in heavy use on a daily basis
- FreeBSD has served us well
    - As a computing platform
    - and as a cluster research environment
- FreeBSD is an excellent cluster OS
    - Especially for a FreeBSD shop

**THE AEROSPACE CORPORATION**

# Disclaimer

- All trademarks, service marks, and trade names are the property of their respective owners.