# International Outlook

Maurizio Davini

Dipartimento di Fisica e INFN Pisa

*maurizio.davini@df.unipi.it*

# Components

- Extensions for clustering
- Configurations,installation and administration
- Monitoring
- Batch e schedulers
- Cluster File Systems

# Extension for clustering

Workshop CCR INFN, La Biodol

# Bproc

- http://bproc.sourceforge.net/
- BProc, The Beowulf Distributed Process Space, is a set of kernel modifications, utilities, and libraries which allow a user to start processes on other machines in a Beowulf-style cluster. Remote processes started with this mechanism appear in the process table of the front end machine in the cluster.

# OpenMosix

- http://www.openmosix.org/
- is a software package that enhances the Linux kernel with cluster capabilities. MOSIX allows for the automatic and transparent migration of processes to other nodes in the cluster, while standard Linux process control utilities, such as 'ps' will show all processes as if they are running on the node the process originated from

# SSI

◆ http://ssic-linux.sourceforge.net/

◆ Single System Image Clusters for Linux (SSI) aims at providing a full, highly available, single system image cluster environment for Linux, with the goals of availability, scalablity, and manageability, built from standard servers.

# CI (cluster Infrastructure)

- http://ci-linux.sourceforge.net/
- Cluster Infrastructure for Linux (CI) aims at developing a common infrastructure for Linux clustering by extending cluster membership and internode communication subsystems from Compaq's NonStop Clusters for Unixware code base. This project also provides the basis for the SSI Clusters for Linux project.

# Configuration and Management

# CFEngine

◆ http://www.gnu.org/software/cfengine/cf

◆ Cfengine is a tool for administration and configuration of large (or small) networks of computers. It uses the idea of classes to define and automate the configuration and maintenance of systems.

# LCFG

◆ http://www.lcfg.org/

◆ LCFG (Local Configuration System) is a system for automatically installing and managing the configuration of large numbers of Unix systems. It is particularly well suited to environments with diverse and rapidly changing configurations

# C3

- Cluster Control e Command Tool
- http://www.csm.ornl.gov/torc/C3/
- The C3 (Cluster Command and Control) tool suite was developed at Oak Ridge National Lab and implements a number of command line based tools that aid in the operation and management of clusters.

# System Installation Suite

- http://sisuite.org/
- The System Installation Suite is the grand project name for the joint effort of System Configurator, System Imager, and System Installer. All three are being developed to work together to provide an easy to use solution for installing and configuring a cluster of heterogenous hardware and to be distribution agnostic

# ImgDisk

◆ http://www.unix.mcs.anl.gov/systen

◆ Imgdisk can save and restore disk images: the partition table, MBR, partition boot records, ext2 file system contents, and can also restore swap partitions.
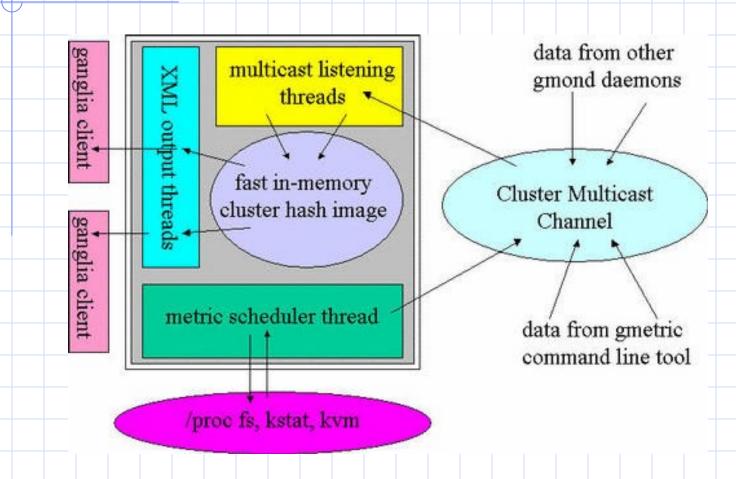
# Webmin

◆ http://www.webmin.com/webmin/

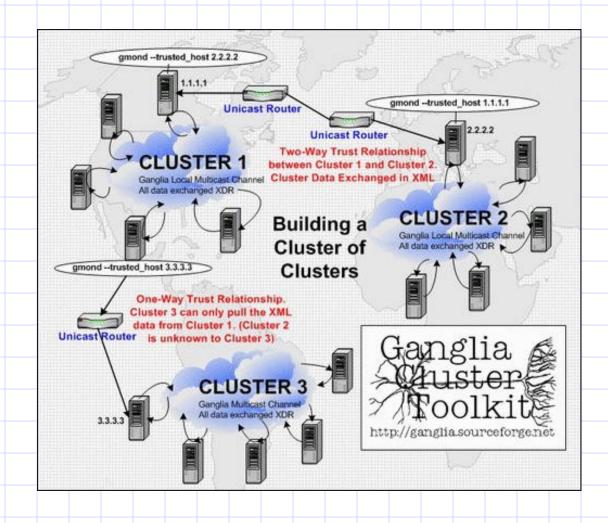◆ Webmin is a web based interface for system administration of Unix systems.

# Monitoring

# GANGLIA

- http://ganglia.sourceforge.net/
- The Ganglia Cluster Toolkit is a massively scalable cluster monitoring and execution environment. Currently, after the monitoring system has been released the execution environment too.

# Ganglia

# Ganglia

# SuperMon

◆ http://www.acl.lanl.gov/supermon/

◆ Supermon is a high speed monitoring system for large scale clusters

◆ Recent benchmarks have shown that supermon can achieve single node monitoring speeds that are significantly higher than previous methods - we have observed 6000 samples-per-second with supermon, while the same benchmark was only able to achieve 275 samples-per-second with older methods used by RPC.rstatd.

# Nagios

- http://www.nagios.org/
- The Nagios Network Monitor can monitor network services and host resources, contains contact notifications when problems occur (via email, pager, or user-defined method), contains an optional web interface for viewing network status, notification, problem history or log files, and is highly extensible, making it easy for the creation of user-developed service checks

# VACM

- http://www.valinux.com/software/vacm

- VACM (VA Cluster Manager) is a tool for monitoring and managing a large cluster of nodes. There is also a page for VACM on SourceForge, located here.

# Clunix

- http://www.cmap.polytechnique.fr/~sylva

- The clunix tools are a free, GPL'd set of utilities for managing a Linux cluster. The components are a perl based load balancer, a php based web interface for monitoring the state of the cluster, and a daemon.

# ECT

- http://www.alphaworks.ibm.com/tech/ect
- Enhanced Cluster Tools for Linux (ECT) is a set of tools which complement IBM's Cluster Systems Management (CSM) and enhance the management of clusters, providing features such as remote access to hardware inventory and vitals, and remote access to service processor logs.

# Cluster File Systems

# ClusterNFS

- http://clusternfs.sourceforge.net/
- ClusterNFS is a set of patches for the UNFSD server to allow multiple diskless clients to mount the same root filesystem, by providing "interpreted" filenames.

# GFS   $$

◆ http://www.sistina.com/products_gfs

◆ GFS (Global File System) is fault tolerant and distributed. It is also both a cluster and a journalling file system.

# GPFS    $$

◆ IBM's General Parallel File System (GPFS) allows users shared access to files that may span multiple disk drives on multiple nodes. It offers many of the standard UNIX ® file system interfaces allowing most applications to execute without modification or recompiling. UNIX file system utilities are also supported by GPFS.

# PVFS

- http://parlweb.parl.clemson.edu/pvf
- PVFS (Parallel Virtual File System) is being developed at Clemson University's Parallel Architecture Research Lab and is closely tied to the Beowulf (currently a.k.a. Scyld) Project.

# LUSTRE

- http://www.lustre.org/
- Lustre is a novel storage and file system architecture that aims at building a next-generation cluster file system, to service clusters with 10,000s of nodes, petabytes of storage, and move 100s of GB/sec, as well as offering security and management.

# Batch systems e schedulers

# Sun Grid Engine

- http://gridengine.sunsource.net/
- The Gride Engine Project is an open source project based on Sun's commercial product, "Sun Grid Engine," which can be seen here. Grid Engine is Distributed Resource Management software, used to create compute farms.

# LSF $$

- http://www.platform.com/
- LSF (Load Sharing Facility) is a suite of software available for various Unixes and NT. It performs load sharing and balancing, and job scheduling.

# PBS $$

- http://pbs.mrj.com/
- PBS (Portable Batch System) is a batch queueing and load balancing system originally developed for NASA. It is available for a variety of Unix platforms. There is an Open Source version of PBS, called OpenPBS, which is located here.

# Condor

- http://www.cs.wisc.edu/condor
- Condor is a software package that does job scheduling and load balancing. It is available for most Unixes, and a port to NT is currently underway.

# MAUI

◆ http://supercluster.org/maui/

◆ Maui is an advanced scheduler for clusters and supercomputers. It can support various fairness policies, dynamic priorities, extensive reservations, and fairshare

# Installation systems

# Cluster distributions

- http://clusters.top500.
- Why are 26% rolling there own clusters?
  - This is the real message of the poll

**Poll Results**

**What Cluster system(Distribution) do you use?**

Other 26%
Oscar 23%
Score 13%
Scyld 12%
MSC.Linux 12%
NPACI Rocks 8%
SCE 6%

302 votes | 3 comments

# LTSP

- http://www.ltsp.org/
- The name pretty much says it all. The LTSP is a way to run thin client computers in a Linux environment.

- Clustering openMosix as an option

# FAI

- http://www.informatik.uni-koeln.de/fai
- FAI (Fully Automatic Installation) is a non-interactive system to install Debian GNU/Linux on a cluster. It is a collection of Perl and shell scripts, and will work with a variety of PC hardware

# KA

- http://ka-tools.sourceforge.net/
- Ka is a toolkit designed to install and administer a cluster of boxes. It focus on scalability of parallel system installation, data distribution and process launching. Ka has been tested on clusters up to 225 nodes

# Clubmask

- clubmask.sourceforge.net
- Clubmask is a collection of existing Open Source and new software for the installation, configuration, and management of Beowulf style high performance computing clusters. The design and goal of the project is to provide a "Physicist Proof", completely turnkey set of tools

# Clubmask

1. Bproc - unified cluster process namespace and control
2. Cfengine - class based configuration tool
3. Kickstart - RedHat automated installation scripts
4. Lam - MPI environment
5. Maui Scheduler - advanced batch scheduler with a large feature set well suited for high performance computing (HPC) platform

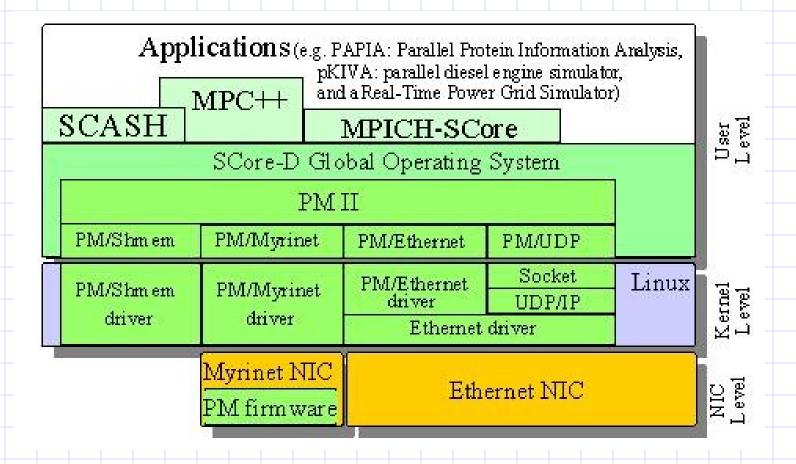6. ZODB - Embedded Python Object

# ClusterIt

- http://www.garbled.net/clusterit.html
- ClusterIt is a collection of clustering tools, allowing for heterogeneous cluster makeup, including various architectures and operating systems, various authentication mechanisms, job sequencing, distributed virtual terminals, and more.

# SCORE

◆ http://www.pccluster.org/

◆ SCore, by Real World Computing Partnership (RWCP) is not a Beowulf style cluster in the sense that SCore software is designed for the high performance cluster environment without using the TCP/IP stack.

# Score

# WareWulf

◆ http://www.runlevelzero.net/greg/warewu

◆ Warewulf is a distribution of tools that are designed to aid in the implementation of Beowulf style clusters. The software is a bootable ISO image that is easily modified, and slave node filesystems can be booted off the CDROM, so no changes must be made to the existing hard disk if wanted.
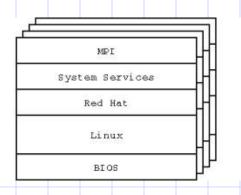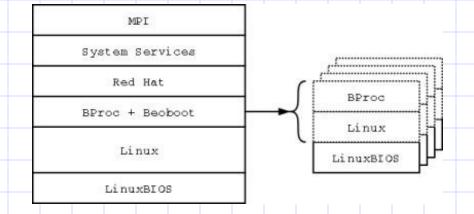
# Cplant

- http://www.cs.sandia.gov/cplant/
- Computational Plant (a.k.a. Cplant) is a newly released project coming from the folks at the Sandia National Laboratories. The goal is "to provide a commodity-based, large-scale computing resource that meets the level of compute performance needed by Sandia's critical applications."

# ClusterMatic

◆ http://www.clustermatic.org/

◆ Clustermatic is a collection of technologies being developed at the Cluster Research Lab at Los Alamos National Laboratory. Besides the new software being developed by the group, existing projects such as LinuxBIOS and BProc are integrated into it as well.

# ClusterMatic

# OSCAR

- http://oscar.sourceforge.net/
- OSCAR (Open Source Cluster Application Resources) is a bundle of software designed to make it easy to build, maintain, and use a Linux cluster

# OSCAR – An Overview

- Open Source Cluster Application Resources
- Cluster on a CD – automates cluster install process
- IBM, Intel, NCSA, ORNL, MSC Software, Dell
- NCSA "Cluster in a BOX" base
- Wizard driven
- Nodes are built over network
- OSCAR <= 64 node clusters for initial target
- Works on PC commodity components
- RedHat based (for now)
- Components:  Open source and BSD style license

# Why OSCAR?

- ◆ NCSA wanted "Cluster-in-a-Box" Distribution
  - ▪ NCSA's "X-in-a-Box" projects could lie on top
  - ▪ X = Grid, Display Wall, Access Grid
- ◆ Easier, faster deployment
- ◆ Consistency among clusters
- ◆ Other organizations had the same interest
  - ▪ Intel, ORNL, Dell, IBM, etc.
  - ▪ NCSA jumps on board to contribute to OSCAR

# OSCAR Basics

◆ What does it do?
- OSCAR is a cluster packaging utility
- Automatically configures software components
- Reduces time to build a cluster
- Reduces need for expertise
- Reduces chance of incorrect software configuration
- Increases consistency from one cluster to the next

◆ What *will* it do in the future?
- Maintain cluster information database
- Work as an interface not just for installation, but also for maintenance
- Accelerate software package integration into clusters

# OSCAR Basics

**How does it work?**

◆ version 1.0, 1.1
- ▪ LUI executes
- ▪ = Linux Utility for cluster Install
  - ◆ Network boots nodes via PXE or floppy
  - ◆ Nodes install themselves from rpms over NFS from the server
  - ◆ Post installation configuration of nodes and server executes

◆ version 1.2+
- ▪ SIS = System Installation Suite
  - ◆ System Imager + LUI = SIS
  - ◆ Creates image of node filesystem locally on server
  - ◆ Network boots nodes via PXE or floppy
  - ◆ Nodes synchronize themselves with server via rsycn
  - ◆ Post installation configuration of nodes and server

# Components

◆ OSCAR includes (currently):
- C3 – Cluster Management Tools (ORNL)
- SIS – Network OS Installer (IBM)
- MPI-CH – Message Passing Interface
- OpenSSH/OpenSSL – Secure Transactions
- PBS – Job Queuing System
- PVM – Parallel Virtual Machine

◆ Current Prerequisites:
- Networked PC hardware with disk drives
- Server machine with Redhat installed
- Redhat CD(s)
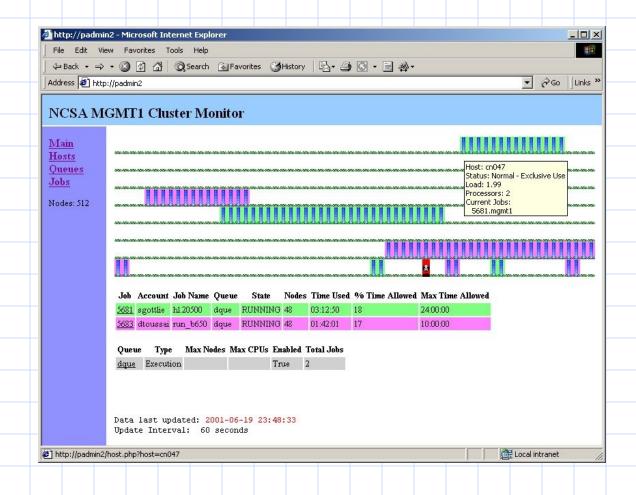- 1 head node + N compute nodes

# Installation Overview

- Install RedHat
- Download OSCAR
- Print/Read document
- Copy RPMS to server
- Run wizard (install_cluster)
  - Build image per client type (partition layout, HD type)
  - Define clients (network info, image binding)
  - Setup networking (collect MAC addresses, configure DHCP, build boot floppy)
  - Boot clients / build
  - Complete setup (post install)
  - Install test suite
- Use cluster



OSCAR Installation Wizard

Welcome to the OSCAR wizard!

Step 1: Build OSCAR Client Image    Help
Step 2: Define OSCAR Clients    Help
Step 3: Setup Networking    Help

Before continuing, network boot all of your nodes. Once they have completed installation, reboot them from the hard drive. Once all the machines and their ethernet adaptors are up, move on to the next step.

Step 4: Complete Cluster Setup    Help
Step 5: Test Cluster Setup    Help

Quit

# OSCAR 2

- Major Changes - Summary
  - No longer bound to OS installer
  - Components are package based, modular
  - Core set of components mandatory
  - API established and published for new packages
  - Package creation open to community
  - Database maintained for node and package information
  - Add/Remove Node process will be improved
  - Web based wizard
  - Scalability enhancements
  - Security Options
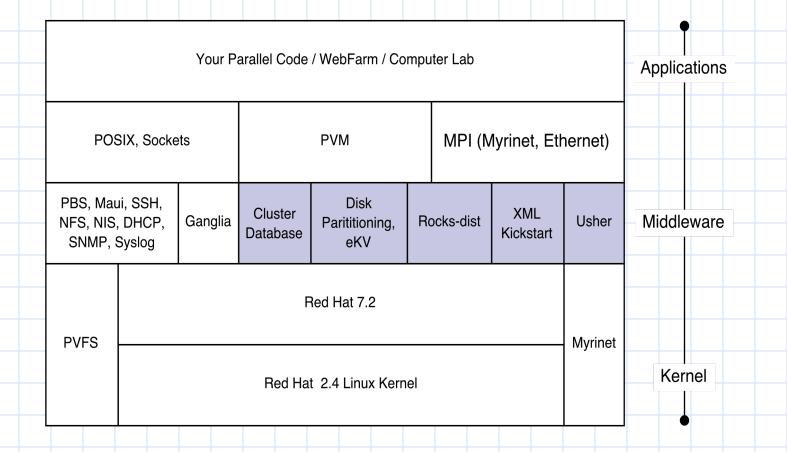  - Support more distributions and architectures
  - New Features

# CluMon

# ROCKS

◆ http://rocks.npaci.edu/

◆ The Rocks Clustering Toolkit, from the folks at NPACI, is a collection of Open Source tools to help build, manage, and monitor, clusters.

# Who is Using It?

◆ Growing (and partial) list of users that we know about:

- SDSC, SIO, UCSD (8 Clusters, including CMS (GriPhyN) prototype)
- Caltech
- Burnham Cancer Institute
- PNNL (several clusters, small, medium, large)
- University of Texas
- University of North Texas
- Northwestern University
- University of Hong Kong
- Compaq (Working relationship with their Intel Standard Servers Group)
- Singapore Bioinformatics Institute
- Myricom (Their internal development cluster)
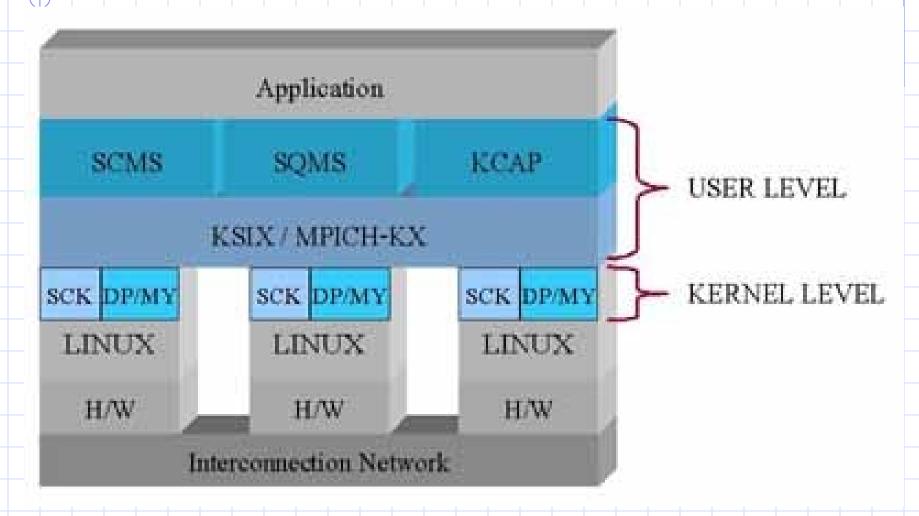
# Major Components

| Your Parallel Code / WebFarm / Computer Lab | | | | | | | Applications |
|---|---|---|---|---|---|---|---|
| POSIX, Sockets | | PVM | | MPI (Myrinet, Ethernet) | | | |
| PBS, Maui, SSH, NFS, NIS, DHCP, SNMP, Syslog | Ganglia | Cluster Database | Disk Parititioning, eKV | Rocks-dist | XML Kickstart | Usher | Middleware |
| PVFS | Red Hat 7.2 | | | | | Myrinet | |
| | Red Hat 2.4 Linux Kernel | | | | | | Kernel |

# SCE

- http://www.opensce.org/
- SCE (Scalable Cluster Environment) is an easy to use set of interoperable Open Source tools that allows the user to quickly install, configure, and use, a Beowulf cluster.

# SCE

- AMATA High Availability support for Beowulf Cluster
- Beowulf Diskless Cluster Suite a utility that help you build a diskless Beowulf cluster easily.
- KCAP Web and VRML based system monitoring and navigation tool for large scale cluster
- KSIX Middle ware layer that offer powerful programing extension in cluster environment
- SCMS SMILE Cluster Management System . A powerful system administration tools for Beowulf style cluster
- SQMS Flexible and extensible batch scheduling system for beowulf

# SCE

# MSC.Linux

- http://www.msclinux.com/software/
- MSC.Linux is a high performance/cluster distribution that is designed for computational environments in engineering and life sciences.

# ClusterWorX **$$**

- http://www.linuxnetworx.com/produc...
php
- ClusterWorX, from the folks at Linux NetworX, is a management and monitoring package, that supports customized monitoring and notification, disk cloning, node management, and more

# Alinka Raisin **$$**

◆ www.alinka.com

◆ Alinka's Raisin software package can do everything from creation to administration of High Performance Linux clusters. It uses Batch Queuing systems such as PBS (or others such as LSF or NQE upon request), Mosix for process migration, and the Parallel file systm PVFS, and comes with a web based user interface

# Cluster System Management(CSM)

- http://www1.ibm.com/servers/eserver/clusters/sc
- Cluster Systems Management for Linux (CSM) enables organizations to build, manage, and expand clusters. A free demonstration of CSM, called the Cluster Starter Kit, is available here.

# SCALI **$$**

◆ http://www.scali.com/products/ssp.html

◆ The Scali Software Platform (SSP) delivers a number of tools for ease of installation, administration, maintenance, and operational use of clusters ranging from a handful to hundreds of nodes, that targets all aspects of building, maintaining, and using a cluster. It covers everything from low level drivers to high level administration.

# SCYLD Beowulf **$$**

- http://www.scyld.com/
- The Scyld Beowulf Cluster Operating System is the second generation of The Beowulf clustering software. Scyld Computing Corporation was started by Donald Becker and few other folks from the original Beowulf Project team.

# Qlusters      $$

- www.qlusters.com
- For the openMosix Enthusiasts....

# Qlusters OS features (1)

- Based in part on openMosix technology
- Migrating sockets
- Network RAM already implemented
- Cluster parallel Installer,
- Cluster Configurator,
- Qsense ( automatic detection of nodes no-more /etc/mosix.map )
- Cluster Monitor (written in Flash)
- Logical Clusters ( separate domain of clustering and administration)

# Qlusters OS features (2)

◆ Automatic update

◆ New Load Balancer

◆ Threaded applications migration

◆ Linux kernel 2.4.18 with (VM by A.Arcangeli integrated with Reverve Mapping by R.V.Riel)

◆ Over 100 patches ( RedHat Quality)

◆ Kernel latency reduced by 65% due to Robert Love latest pre-emption patch

# Qlusters OS features (3)

- Queue Manager ,Launcher, Scheduler
- Job Description Language in XML
- Integration with GFS completed
- Integration with AFS planned
- IBM xSeries NUMA support
- DSM in a few months
- Port to IA64 underway

# Qluster Os features (4)

- grid with multiplatform consideration (recompiles when transferring on a cluster of different architecture )

- Grid component (link Q-OS clusters),OS independent (AIX,Solaris…)
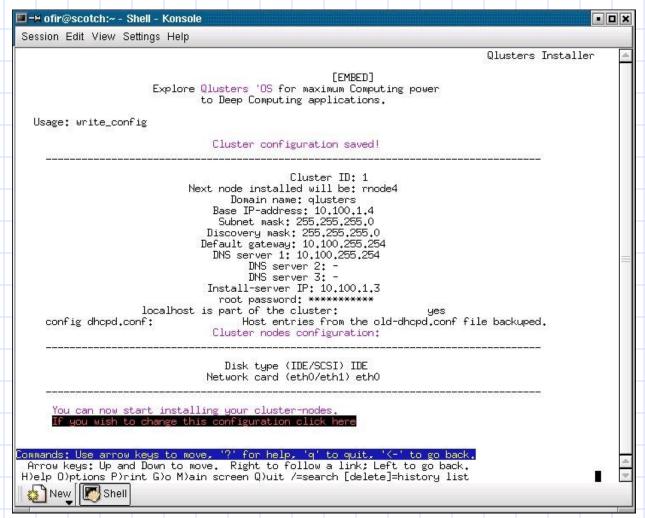
- **Will be the RedHat HPC cluster solution**?

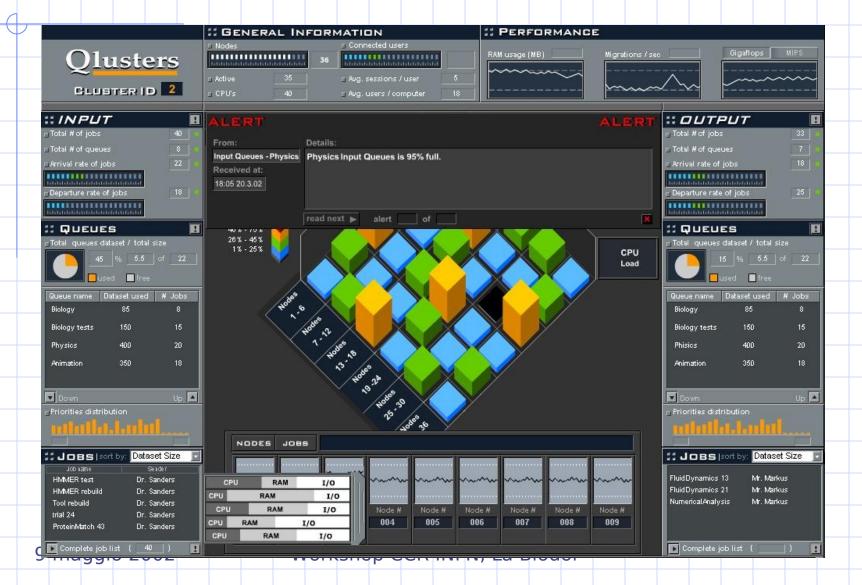# The Installer

# The Installer

# The Installer

# The Installer

# The Monitor

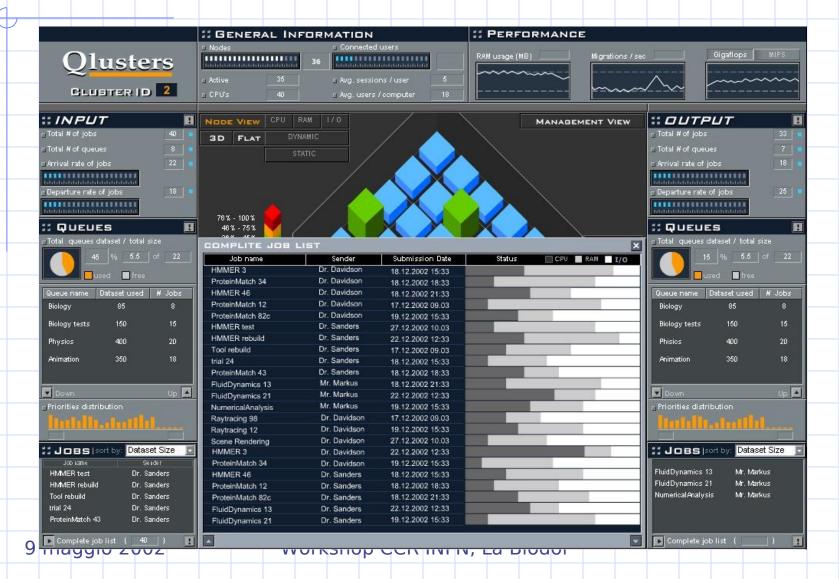# The Monitor

# Interesting sites

- www.lcic.org
- www.beowulf-underground.org
- www.alinka.com ( per Alinka Clustering Newsletter)

# The End