# Scientific Workflow Management

## CCSM Software Engineering Working Group Session

### 6/19/2008

### Scott Klasky

R. Barreto, C. Jin, J. Lofstead, M. Parashar,
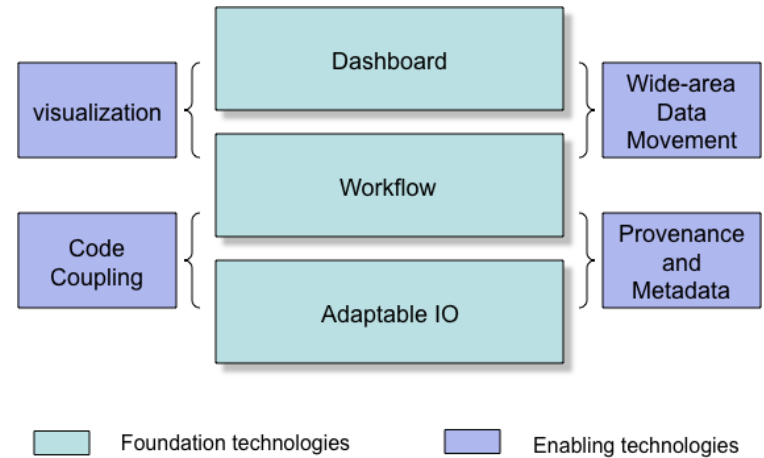N. Podhorszki, K. Schwan, A. Shoshani,
M. Vouk, M. Wolf

# Outline

- ## ADIOS.

- ## Workflow.
  - What is a workflow
  - What advantages over python.
  - Monitoring workflow.
  - Coupling workflow.
  - Movie.
  - Marriage of ADIOS + workflow.
  - Napkin drawing of climate workflow.

- ## Dashboard.

- ## Conclusions.

# End to End Computing at ORNL

- Combines
  - Petascale Applications.
  - Petascale I/O techniques.
  - Workflow Automation.
  - Provenance capturing system.
  - Dashboards for real-time monitoring/controlling of simulations, and creating social spaces for scientists.
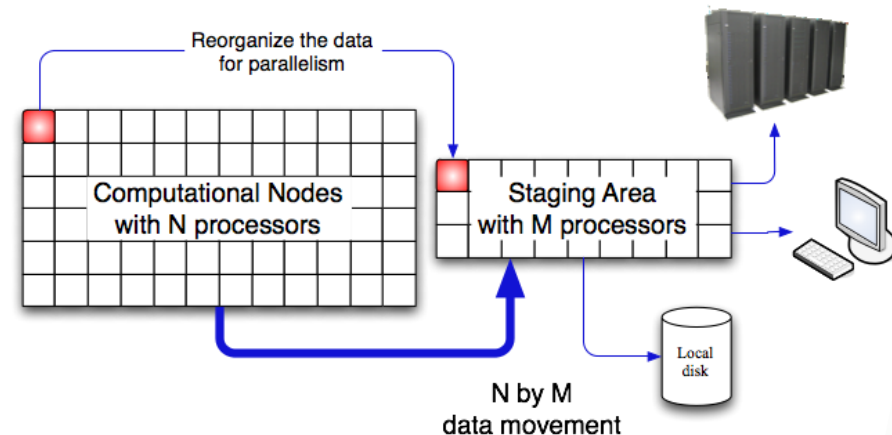


- **Approach**: place highly annotated, fast, easy-to-use I/O methods in the code, which can be monitored and controlled, have a workflow engine record all of the information, visualize this on a dashboard, move desired data to user's site, and have everything reported to a database.

# ADIOS Overview – Design Goals

- ADIOS is an I/O componentization, which allows us to
  - Abstract the API from the IO implementation.
  - Switch from synchronous to asynchronous I/O at runtime.
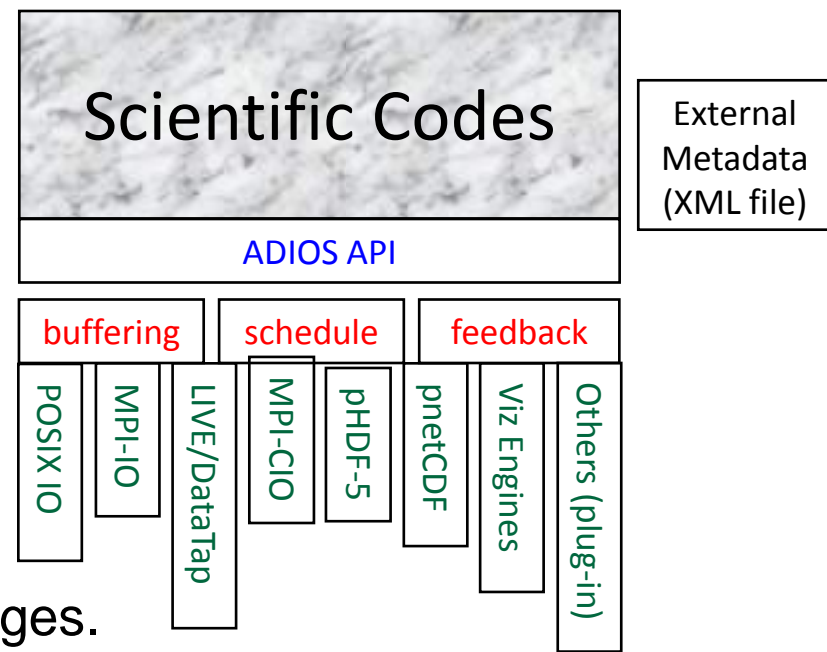  - Change from real-time visualization to fast I/O at runtime.

- Combines.
  - **Fast** I/O routines.
  - **Easy** to use.
  - **Scalable** architecture (100s cores) millions of procs.
  - **QoS.**
  - Metadata rich output.
  - Visualization applied during simulations.
  - Analysis, compression techniques applied during simulations.
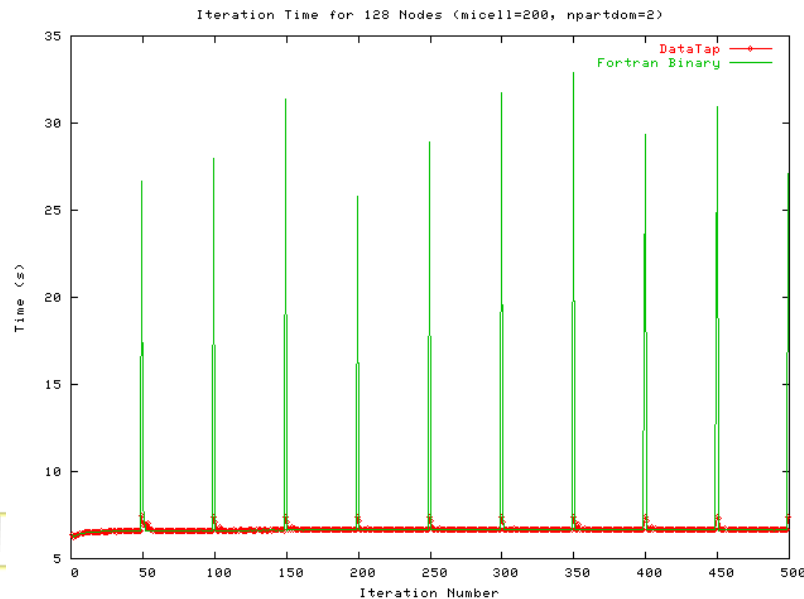  - Provenance tracking.

# ADIOS Overview

- Overview
  - Allows plug-ins for different I/O implementations.
  - Abstracts the API from the method used for I/O.

- Simple API, almost as easy as F90 write statement.

- Both synchronous and asynchronous transports supported without code changes.

- Best practices/optimize IO routines for all supported transports "for free"

- Componentization.
  - Don't worry about IO implementation.
  - Components for IO transport methods, buffering, scheduling, and eventually feedback mechanisms.

- Change I/O method by changing XML file only, with metadata inside.

- Will support strong-coupling (code-code) in the future using a shared global memory address space

**Scientific Codes**

External Metadata (XML file)

ADIOS API

buffering | schedule | feedback

POSIX IO | MPI-IO | LIVE/DataTap | MPI-CIO | pHDF-5 | pnetCDF | Viz Engines | Others (plug-in)

GPSC   SDM CENTER   Georgia   RUTGERS   CPES Center for Plasma Edge Simulation   GSEP   NC State   Sandia National Laboratories   NORTHWESTERN UNIVERSITY   Office of Science U.S. DEPARTMENT OF ENERGY   OAK RIDGE National Laboratory

# Initial ADIOS performance.

- ## MPI-I/O method.
  - GTC and GTS codes have achieved over 25 GB/sec on Cray XT at ORNL.
    - 30GB diagnostic files every 3 minutes, 1.2 TB restart files every 30 minutes, 300MB other diagnostic files every 3 minutes.
  - Chimera code speed up by 6.5% (overall time).

- ## DART: <2% overhead for writing 2 TB/hour with XGC code.

- ## DataTap vs. Posix
  - 1 file per process (Posix).
  - 5 secs for GTC computation.
  - ~25 seconds for Posix I/O
  - ~4 seconds with DataTap



Iteration Time for 128 Nodes (micell=200, npartdom=2)

GPSC

SDM CEN

Georgia

RUTGERS

$\ell$ PES
Center for Plasma Edge Simulation

National Laboratory

# EFFIS is a marriage between Keplerand ADIOS and the dashboard.

- ADIOS is being modified to send the I/O (and coupling) metadata from the compute nodes over to Kepler.
  - We send the file information( including the path).
  - We send the metadata that is contained.
    - Variables, + other annotations.
  - We can send commands to inform Kepler what to do.

- The information is NOT sent if there is no one listening; i.e. not dependent on Kepler).

- The information is saved in a database.

# What is the Kepler Workflow Framework?

**Kepler is a proven DOE technology from the SDM center for orchestrating scientific workflows, which aid construction and automation of scientific problem-solving processes.**

- **Kepler workflow framework**
  - **Captures provenance information for**
    - Data provenance (Where did my data come from?)
    - Data movement and data replication (*e.g.*, during code coupling)
    - Tar files stored on HPSS (at NERSC or ORNL)
    - Workflow actions saved in log files for user debugging
  - **Is more powerful than Python scripts**
    - Allows pipeline-parallel processing with ease
    - Allows work to continue even if some scripts/components fail
    - Allows checkpoint/restart of the workflow
    - Easy to modify workflow for a continuously changing group of scientists
  - **Provides an excellent connection to databases**
    - Allows for easy queries of shots from coupled simulations
    - Large SDM effort to save provenance data into database

GPSC  SDM CENTER  Georgia  RUTGERS  CPES Center for Plasma Edge Simulation  GSEP  NC State  Sandia National Laboratories  NORTHWESTERN UNIVERSITY  Office of Science U.S. DEPARTMENT OF ENERGY  OAK RIDGE National Laboratory

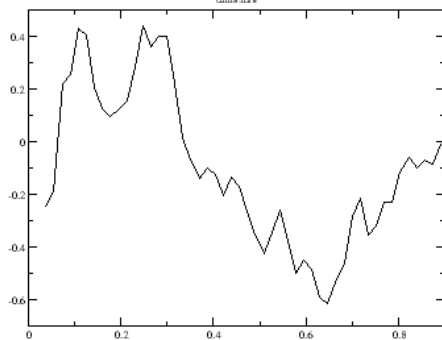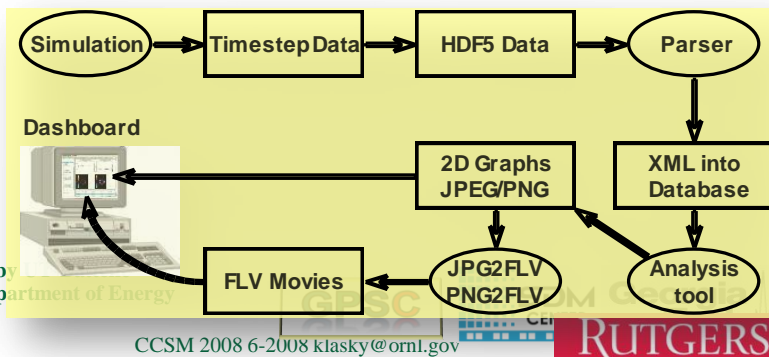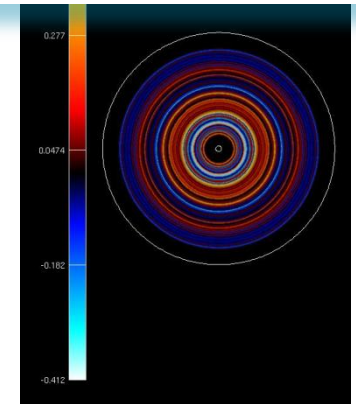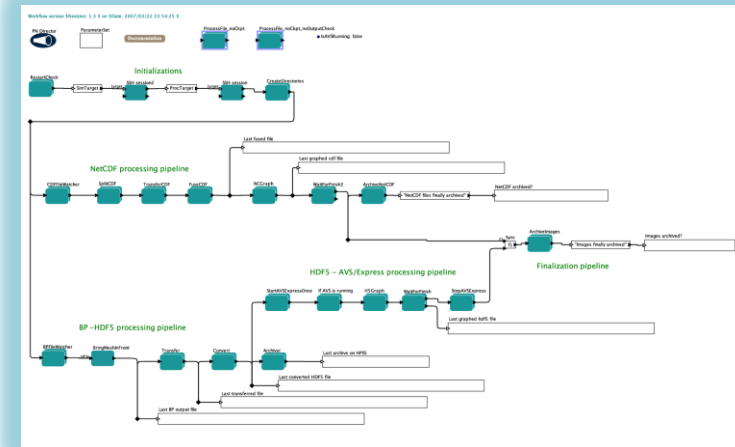# Workflow automation needs in Center for Plasma Edge Simulation



- Automate the data processing pipeline and simulation monitoring
  - Transfer of simulation output to remote machine
  - Execution of conversion routines
  - Image creation, data archiving, dynamic monitoring
- Automate the code coupling pipeline
  - Run simulation on a large supercomputer
  - Check linear stability on another machine
  - Re-run simulation if needed
- Requirements for petascale computing
  - Easy to use
  - Parallel processing
  - Dashboard front-end
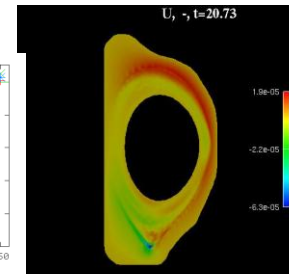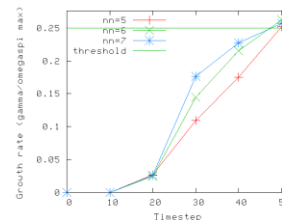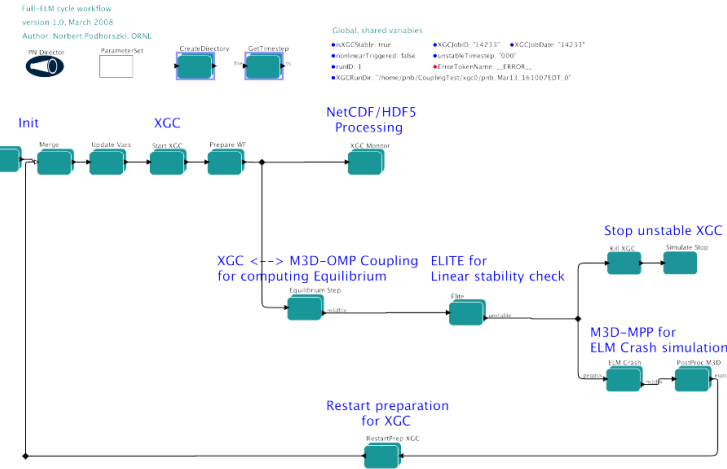  - Robustness
  - Dynamic monitoring
  - Configurability

# Workflows for monitoring a simulation

- NetCDF files
  - Transfer files to e2e system on-the-fly
  - Generateimages using grace library
  - ArchiveNetCDF files at the end of simulation

- Binary files from ADIOS
  - Transfer to e2e system using *bbcp*
  - Convert to HDF5 format
  - Generateimages with AVS/Express (running as service)
  - Archive HDF5 files in large chunks to HPSS

- Record Provenance information for everything!

# Coupling Fusion codes for Full ELM, multi-cycles



- Run XGC until ELMS are unstable

- M3D coupling data from XGC

  - Transfer to end-to-end system

  - Execute M3D: compute new equilibrium

  - Transfer back the new equilibrium to XGC

  - Execute ELITE: compute growth rate, test linear stability

  - Execute M3D-MPP: to study unstable states (ELM crash)

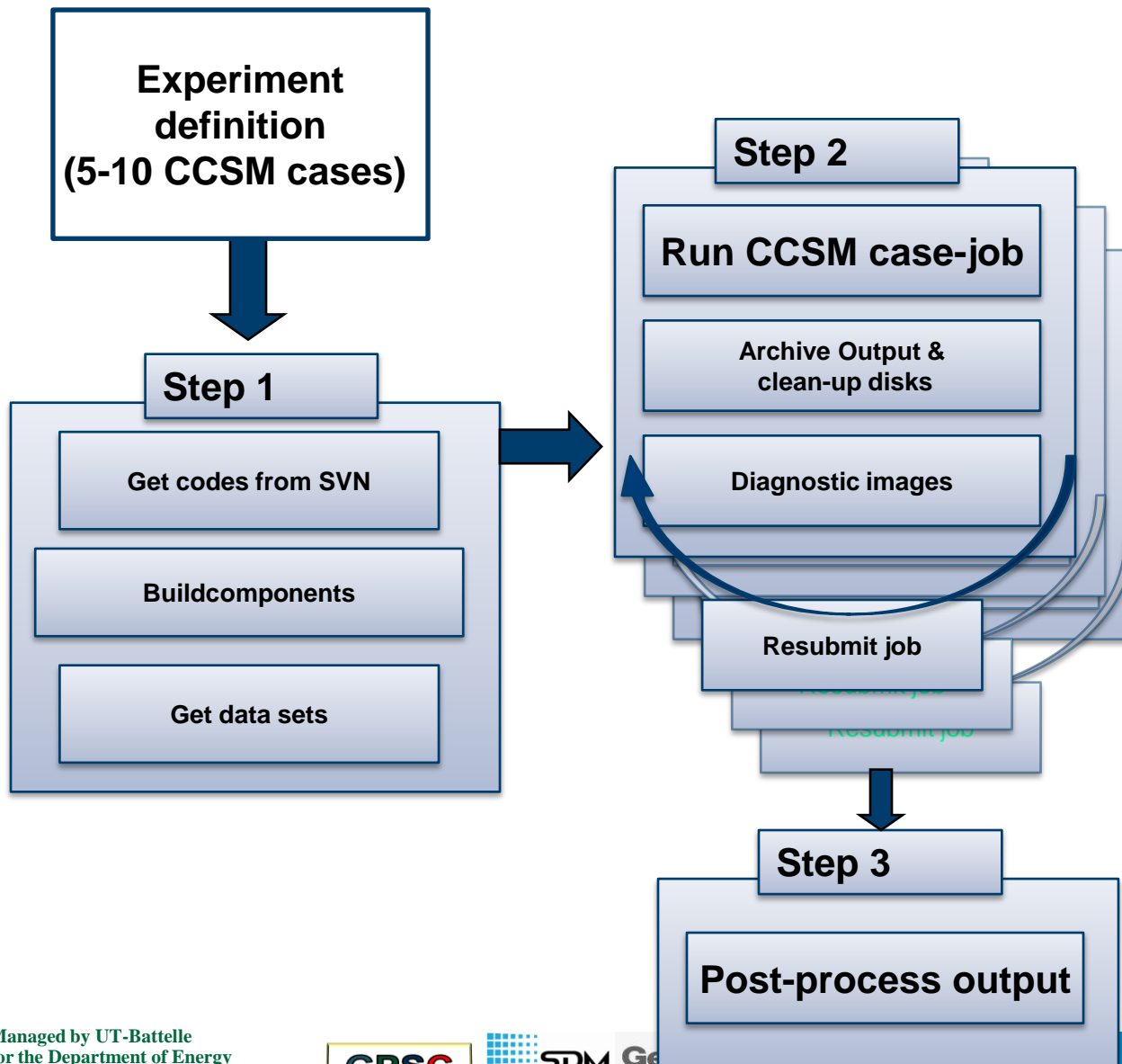  - Restart XGC with new "helaled" equilibrium from M3D-MPP

Jaguar - XGC-0

File   Edit   View   Terminal   Tabs   Help

| Jaguar - XGC-0 | Ewok - M3D and Worfklow processing | Local - Workflow GUI |

pnorbert@jaguar12:/tmp/work/pnorbert/fullelm>

**XGC is submitted on JAGUAR.**

# Climate workflow

**Experiment definition (5-10 CCSM cases)**

**Step 1**
- Get codes from SVN
- Buildcomponents
- Get data sets

**Step 2**
- **Run CCSM case-job**
- Archive Output & clean-up disks
- Diagnostic images
- Resubmit job
- Resubmit job
- Resubmit job

**Step 3**
- **Post-process output**

- Run an ensemble of several cases at once
  - ✓ one case consists of many jobs executed repetitively
- Run jobs at NCCS and other resource
  - ✓ one-time-password access from outside
- Archive output on the fly to HPSS while the job is running
- Make images from diagnostic output on the fly and put on the dashboard
  - ✓ to monitor the current status of each case

**Managed by UT-Battelle for the Department of Energy**

CCSM 2008 6-2008 klasky@ornl.gov

GPSC  SDM CENTER  Ge...  RUTGERS  Center for Plasma Edge Simulation  CCELP  Sandia National Laboratories  NORTHWESTERN UNIVERSITY  Office of Science U.S. DEPARTMENT OF ENERGY  OAK RIDGE National Laboratory

# Design Criteria for the Dashboard

- **Goal:** provide users an easy way over the web to dynamically monitor simulation progress, to view images and movies, to perform basic analysis, and to move files to their site

- New design criteria for FSP codes on leadership class computers
  - Must support very large and small data, in a scalable fashion

- New security with One Time Passwords
  - Unrealistic to think that we can monitor jobs via one type of data output
  - Unrealistic to think that we can move data from a large parallel disk to user space

- Data management must be incorporated into the design
  - Database back-end is as important as front-end
  - Provenance display is very important to monitor long-running jobs

- Must be able to monitor computers/jobs from all resources

- Need to plug-in new visualization routines into the display

- Need to plug-in new analysis routines into the system

- Need to collaborate via shared space

- Make it robust by using enterprise web-2 technologies

# Dashboard: Simulation Monitoring

- A basic browser-based visualization tool

- Local interaction for graphs/visualizations

- Server used for data manipulation, extraction

- Server used for more complex analysis

- Allows for pan and zoom locally

- Asynchronously queries

- MySQL for database support

# Dashboard: Job Monitoring

Monitors machines, simulations and DB

- Secure login with OTP

- Job submission and kill

- Search old jobs

- See collaborators jobs

- Annotations/Notes

- Text display/movies

# Dashboard movie

# Collaborative Analysis Features

- Basic analysis on dashboard will feature
  - Calculator for simple math, done in Python
  - Hooks into "R" for pre-set functions
  - Ability to save the analysis into a new function
  - 2d and time history plots (initial version)
  - Full 3d plots (in future version)

- Advanced analysis will contain
  - Parallel backend to VisIT server, VisTrails, Parallel R, and custom MPI/C/F90 code
  - We will allow users to place executable code into the dashboard

- In progress: a portable dashboard!

# Conclusions

- ADIOS is an I/O componentization.
  - ADIOS is being integrated integrated into Kepler.
  - Achieved over 20 GB/sec for several codes on Jaguar.
  - Used daily by CPES researchers.
  - Can change IO implementations at runtime.
  - Metadata is contained in XML file.

- Kepler is used daily for
  - Monitoring CPES simulations on Jaguar/Franklin/ewok.
  - Runs with 24 hour jobs, on large number of processors.

- Dashboard uses enterprise (LAMP) technology.