# Performance Evaluation and Analysis Consortium (PEAC) End Station

Presented by

## Patrick H. Worley

Computational Earth Sciences Group
Computer Science and Mathematics Division

SC07

OAK RIDGE
National Laboratory

# Overview

The PEAC End Station provides the performance evaluation and performance tool developer communities access to the Leadership Computing Facility (LCF) systems.

| Consortium goals | |
|---|---|
| System evaluation | • Evaluate the performance of LCF systems using standard and custom micro-, kernel, and application benchmarks |
| Performance tools | • Port performance tools to LCF systems and make them available to National Center for Computational Sciences (NCCS) users<br><br>• Further develop the tools to take into account the scale and unique features of LCF systems |
| Performance modeling | • Validate the effectiveness of performance modeling methodologies<br><br>• Modify methodologies as necessary to improve their utility for predicting resource requirements for production runs on LCF systems |

OAK RIDGE
National Laboratory

# Overview (continued)

| Consortium goals (continued) | |
| --- | --- |
| Application analysis and optimization | • Analyze performance<br>• Help optimize current and candidate LCF application codes |
| Performance and application community support | • Provide access to other performance researchers who are interested in contributing to the performance evaluation of the LCF systems or in porting complementary performance tools of use to the NCCS user community<br>• Provide access to application developers who wish to evaluate the performance of their codes on LCF systems |

## All of this must be accomplished while adhering to the "Golden Rules" of the performance community:

- Low visibility (no production runs!)
- Open and fair evaluations
- Timely reporting of results

OAK RIDGE
National Laboratory

# Status as of 8/28/07

**32 active users,
39 active projects:**

- 13 application performance analysis and optimization

- 8 system evaluation

- 8 tool development

- 6 infrastructure development

- 4 application modeling

**Consuming:**

- XT4: 1,168,000 processor hours (exceeding 1,000,000 processor-hour allocation)

**Contributing to:**

- 1 refereed journal paper

- 1 invited journal paper

- 6 refereed proceedings papers

- 10 proceedings papers

- 2 book chapters

- Numerous oral presentations

OAK RIDGE
National Laboratory

# System evaluation

| LBNL | Memory, interprocess communication, and I/O benchmarks |
|------|--------------------------------------------------------|
| | APEX-MAP system characterization benchmark |
| | Lattice-Boltzman kernels and mini applications |
| | Application benchmarks from Astrophysics (Cactus), Fluid Dynamics (ELBM3D), High Energy Physics (BeamBeam3D, MILC), Fusion (GTC), Materials Science (PARATEC), AMR Gas Dynamics (HyperCLaw) |
| ORNL | Computation, memory, interprocess comm., and I/O benchmarks |
| | Application benchmarks from Astrophysics (Chimera), Climate (CAM, CLM, FMS, POP), Combustion (S3D), Fusion (AORSA, GTC, GYRO, XGC), Molecular Dynamics (NAMD) |
| SDSC | Subsystem probes for system characterization needed for convolution-based performance modeling |
| Purdue Univ. | Computation, memory, and interprocess comm. benchmarks |
| | Application benchmarks from Chemistry (GAMESS),  High Energy Physics (MILC), Seismic Processing (SEISMIC), Weather (WRF) |

OAK RIDGE
National Laboratory

# Performance tools

| | |
|---|---|
| HPCToolkit | Tool suite for profile-based performance analysis |
| Modeling assertions | Performance model specification and verification framework |
| mpiP | MPI profiling infrastructure |
| PAPI | Performance data collection infrastructure |
| Scalasca | Scalable trace collection and analysis tool |
| SvPablo | Performance analysis system |
| TAU | Performance analysis system |
| MRNet | Scalable performance tool infrastructure |

# Application performance analysis and optimization

| | |
|---|---|
| Chombo | AMR gas dynamics model |
| DeCart | Nuclear code |
| FACETS | Framework application for core-edge transport simulation |
| GADGET | Computational cosmology |
| GTC_s | Shape plasma version of GTC gyrokinetic turbulence code |
| NEWTRNX | Neutron transport code |
| PDNS3D/SBLI | Ab initio aeroacoustic simulations of jet and airfoil flows |
| PFLOTRAN | Subsurface flow model |
| PNEWT | Combustion code |

OAK RIDGE
National Laboratory

# Application code scaling, optimization, and/or performance evaluation

| POLCOMS | Coastal ocean model |
|---------|---------------------|
| S3D | Combustion model |
| TDCC-9d | Nuclear code |
| - | Lattice-Boltzman applications |

# System infrastructure

| | |
|---|---|
| cafc | Co-array Fortran compiler for distributed-memory systems |
| GASNet | Runtime networking layer for UPC and Titanium compilers |
| PETSc | Toolset for numerical solution of PDEs |
| PVFS/Portals | PVFS file system implementation on native Portals interface |
| UPC | Extension of C designed for high-performance computing on large-scale parallel systems |
| - | Reduction-based communication library |

OAK RIDGE
National Laboratory

# Performance modeling

| | |
|---|---|
| PMAC | Genetic algorithm-based modeling of memory-bound computations |
| ORNL | NAS parallel benchmarks; HYCOM ocean code |
| Texas A&M Univ. | GTC fusion code |
| Univ. of Wisconsin | Reusable analytic model for wavefront algorithms, applied to NPB-LU, SWEEP3D, and Chimaera |
| | LogGP model for MPI communication on the XT4 |

OAK
RIDGE
National Laboratory

# Subsystem evaluations



I/O performance characterization (LBL)



Dual vs. single core performance
evaluation using APEX-MAP (LBL)

Ratio of time for all processes sending in halo update
to time for a single sender

| System | 4 neighbors | Periodic | 8 Neighbors | Periodic |
|--------|-------------|----------|-------------|----------|
| BG/L | 2.24 | | 2.01 | |
| BG/L, VN | 1.46 | | 1.81 | |
| XT3 | 7.5 | 8.1 | 9.08 | 9.41 |
| XT4 | 10.7 | 10.7 | 13.0 | 13.7 |
| XT4 SN | 5.47 | 5.56 | 6.73 | 7.06 |

Identifying performance anomalies (ANL)



MPI performance characterization (ORNL)

OAK RIDGE National Laboratory

# Application analyses and benchmarks



Scalability optimizations (ORNL)



Performance sensitivities (SDSC)

Processing of genomes into domain maps: need improved load balancing that takes into account scale-free nature of the graphs.

Porting and optimizing new applications (RENCI/NCSA)

# Tool development

SvPablo source code-correlated performance analysis (RENCI)

mpiP callsite profiling (LLNL/ORNL)

SCALASCA trace-based performance analysis (FZ-Jülich, UTenn)

# Co-principal investigators

| Argonne National Laboratory | Lawrence Berkeley National Laboratory | Lawrence Livermore National Laboratory | Oak Ridge National Laboratory | Rice University | University of California–Berkeley |
|---|---|---|---|---|---|
| William Gropp | David Bailey<br>Leonid Oliker | Bronis de Supinski | Jeffrey Vetter<br>Patrick Worley (PI) | John Mellor-Crummey | Kathy Yelick |

| University of California–San Diego | University of Maryland | University of North Carolina | University of Oregon | University of Tennessee | University of Wisconsin |
|---|---|---|---|---|---|
| Allan Snavely | Jeffrey Hollingsworth | Daniel Reed | Allen Malony | Jack Dongarra | Barton Miller |

# Contact

## Patrick H. Worley

Computational Earth Sciences Group
Computer Science and Mathematics Division
(865) 574-3128
worleyph@ornl.gov

## Barbara Helland

DOE Program Manager
Office of Advanced Scientific Computing Research
DOE Office of Science