# Leadership Computing at ORNL
# Challenges and Opportunities

**Arthur S. Bland, Ann E. Baker, R. Shane Canon, G. Al Geist, Ricky A. Kendall, Jeffrey A. Nichols, Julia C. White, Thomas Zacharia**, *Oak Ridge National Laboratory*

**ABSTRACT:** In May 2004, the National Center for Computational Sciences (NCCS) at Oak Ridge National Laboratory was selected as the Leadership Computing Facility by the U.S. Department of Energy. This paper describes Leadership Computing and the reasons it is important for the United States. It gives a description of the NCCS and its roadmap to expand the capability of the center from 50 TeraFLOPS today to a PetaFLOPS computer over the next three years. Finally, it describes some of the challenges in achieving this goal and our plans for meeting these challenges.

## 1. Why is Leadership Computing Important?

The science of the 21st century demands computational capability well beyond what is available today. These demands cannot be met by simply fielding a computer that is number one on the Top500[1] list. Rather, breakthrough science and engineering requires an architecture that is well suited for scientific applications, a computational environment that ensures effective utilization of that architecture for scientific discovery, a best-in-class communications network and data management infrastructure, and teams of leading experts applying this capability to critical research challenges.

Computational simulation has become an essential tool for research and development. Several reports from the scientific computing community[2–5] indicate a need for sustained computing performance of more than 50 teraflops (TF) in 2006, with >1 petaflops (PF) needed by the end of the decade for the United States to maintain its leadership position in many areas of science, research, and engineering.

Such capabilities will enable: understanding the evolution and consequences of severe events, whether natural or anthropogenic; understanding and intervening in the complex molecular interactions and pathways involved in human disease; and predicting and analyzing the aerodynamic and thermodynamic behavior of aircraft, rocket engines, and space vehicles. Market leadership in many technology-based industries (e.g., automotive, aircraft, pharmaceutical, and chemical) depends on advanced computational simulation to improve existing products, develop new products, and shorten time to market.

The major scientific challenges being addressed by the U.S. Department of Energy Office of Science (DOE-SC) research programs illustrate the possibilities created by advances in scientific supercomputing. Using leadership-class computers, we will be able not only to dramatically extend our exploration of the fundamental processes of nature, but also to advance our ability to predict the behavior of complex natural and engineered systems. Significant results are expected in a variety of application areas, as indicated in Table 1.

1

Table 1. Significant application areas, goals, and needs

| Application | Simulation Need | Computing Need (TF) | Significance |
|---|---|---|---|
| Climate Science | Reliable predictions of climate change at the regional (e.g., state) scale, including biogeochemical feedbacks | > 50 | Provides U.S. policymakers with data to support policy decisions; predicts extreme weather conditions arising from changing climate |
| Magnetic Fusion Energy | Understand and predict plasma turbulence, including self-heating of plasma and heat leakage caused by turbulence | > 50 | Crucial for understanding burning plasmas and quantifying prospects for commercial fusion; underpins U.S. decisions about future international fusion collaborations |
| Computational Biology | Understand the molecular processes that provide the mechanism for cells and organisms to adjust to changing conditions and to take advantage of available energy sources | > 50 | Allows the design of innovative new approaches to bioremediation, energy production, and climate management |
| Computational Chemistry | Reliable prediction of the structures, energetics, and reactions of molecules, especially the large molecules involved in nanoscience, catalysis and energy and environmental science | > 100 | Control of chemical processes responsible for energy produced and pollutants released by engines; prediction and mitigation of spread of pollutants in underground plumes and long-term effects of greenhouse gases and stratospheric ozone depletion; control of processing of high-level radioactive wastes |
| Computational Material Science | Develop systematic approaches, based on sound theoretical models, that integrate across many length and time scales, to understand materials properties | > 100 | Foundation for the design of materials with specified properties from first principles |
| Astrophysics | Realistically simulate the explosion of a core collapse supernova for the first time | >> 100 | Understand the origin of the elements |

In his 2006 State of the Union Address, President George W. Bush announced the American Competitiveness Initiative. The goals of this are:

"To encourage American innovation and strengthen our nation's ability to compete in the global economy. This ambitious strategy will increase Federal investment in critical research, ensure that the United States continues to lead the world in opportunity and innovation, and provide American children with a strong foundation in math and science. The *American Competitiveness Initiative* commits $5.9 billion in FY 2007, and more than $136 billion over 10 years, to increase investments in research and development (R&D), strengthen education, and encourage entrepreneurship and innovation."

In the program plan, the president calls for DOE-SC to build the most powerful civilian supercomputer[6].

In May 2006, the Oak Ridge National Laboratory (ORNL) won a competition among DOE-SC laboratories to host the Leadership Computing Facility (LCF) for science. In this proposal, ORNL described a plan for fielding a series of increasingly powerful computer systems leading to a system with a peak performance of one quadrillion calculations per second, also known as one petaflops (PF). These systems will provide a logical growth path for users to scale applications in numbers of nodes, numbers of CPUs per node, and in the I/O techniques needed to move large amounts of data. As part of the federal budget rollout in February 2006, Dr. Raymond Orbach, director of DOE-SC announced his plan to fulfill the promises of Leadership Computing by committing his office to delivering a 250 TF capability system in 2007, and a 1 PF system in 2008.

## 2. What is Leadership Computing?

In an environment where computing power is doubling every 18-24 months and centers

are constantly installing newer and more powerful systems, what exactly is a Leadership System? The definition comes from the *Federal Plan for High-End Computing*[2]:

> "The goal of such [leadership] systems is to provide computational capability that is at least 100 times greater than what is currently available. A limited set of scientific applications (perhaps 10 per year) would be selected and given substantial access to such systems."
> "Leadership Systems would be designed, procured, and administered to provide computing capabilities to scientific researchers years before equivalent systems become commonplace. They would make possible leading-edge science and engineering research for a select set of challenging, high-payoff, and heretofore unsolvable computational problems."

The LCF at ORNL has deployed leadership systems in its first year of operation and continues on a path to provide leadership computing to the open (unclassified) science and engineering community in the U.S. government, academia, and industry.

## 3. Description of the LCF

The LCF is a component of the National Center for Computational Sciences at ORNL. The NCCS was founded in 1992 as part of the High Performance Computing and Communications Act[7]. DOE's Office of Energy Research (later renamed the Office of Science) held a competition to host a new High Performance Computing Research Center. ORNL won this competition and formed the Center for Computational Sciences (CCS). The CCS acquired and operated a series of computer systems including a KSR-1, several Intel Paragon systems including the XP/S-150 that was the fastest computer in the world at the time of its installation in 1995, an IBM Power3 that was the first 1 TF computer installed in DOE-SC, an IBM Power4 that was the

eighth fastest in the world when installed in 2002, a Compaq AlphaServer-SC, and the world's largest Cray X1 vector system.

In the proposal to create the LCF, the NCCS proposed to double the size of the Cray X1 to 6 TF with a further upgrade to an 18 TF X1E, and to install a 25 TF Cray XT3 system. The capability of these systems would grow to 100 TF and 250 TF, and a plan was presented for achieving 1 PF. The Cray X1E and XT3 were installed and accepted during 2005 and are currently the DOE-SC leadership systems.

While computer centers are often measured by the power of their computers, the LCF is much more that just the collection of computer systems. The users of the LCF are on a steep growth path that will require scaling applications by a factor of 40 in a period of three years; a growth rate of 10 times what would be expected through the typical scaling they might expect from Moore's Law growth of the computers.

DOE-SC instituted the Scientific Discovery through Advanced Computing (SciDAC) program[4] five years ago to develop a new group of scientific application codes for the advanced computer systems being deployed today. Although performance and scalability have been improved, even more must be done for these applications to take full advantage of the leadership systems.

The LCF has a strategy for scaling applications so that they can take advantage of the computers as they are delivered. The strategy is known as the Computational End Station (CES). The CES combines the code development teams from a specific science domain with experts in scaling applications and tuning codes to run on the leadership computer systems. These teams will have very large allocations of time on the computer systems to allow them to perform computational experiments, as well as to continue the development and scaling of their applications.

In addition to computer systems and applications, the LCF has the infrastructure required to support the scientist in using the systems. The LCF has several data analysis systems for pre- and post-processing of data and visualization of the results of calculations. The High Performance Storage System (HPSS) is a multi-petabyte data archive for storing results of simulations. The LCF has deployed an advanced networking infrastructure that links the computers to every major research network in the United States, including: DOE's ESnet and UltraScienceNet, the Internet2 that connects to research universities, NSF's TeraGrid and Cheetah networks, and the National Lambda Rail. This is a scalable networking infrastructure that has the capability of being scaled to as much a 4 terabits per second and can provide circuit switched links directly to individual user locations.

## 4. LCF Roadmap to PetaFLOPS

The LCF plan for providing Leadership Systems to the U.S. science community is an aggressive plan that grows the center from 43 TF peak performance in 2005 to 1000 TF (1 PF) in 2008 with a plan for delivering a sustained petaflops system in 2010 or 2011. Figure 1 below graphically depicts the series of systems.
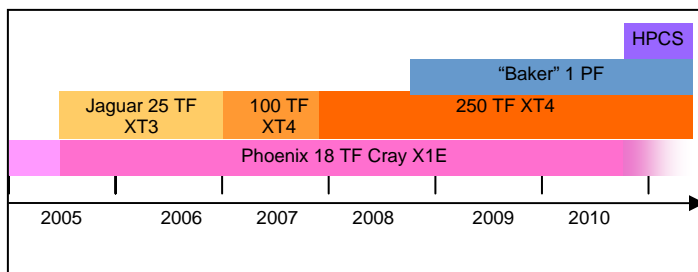


**Figure 1 - LCF Roadmap**

The LCF plan provides a logical sequence of systems that allows the science community to scale their applications through a series of systems, all with the same architecture and programming model. This is a key point

because the work required to scale the applications is substantial without adding in the complexity of porting the applications to new architectures.

The first system installed was the 6 TF Cray X1, called Phoenix. The system is a powerful scalable vector computer that brings back a programming model that immediately benefits a class of applications such as climate models, fusion simulations, and combustion modeling that have high-impact and have not scaled well using typical cluster-based systems. Phoenix was upgraded to a dual-core Cray X1E system in 2005 and now provides 18 TF of capability to the users. Vector processors are a key part of the LCF strategy for achieving a sustained petaflops system so it is important to keep the X1E system viable to allow the researchers to keep their vector applications running.

Also in 2005, the LCF installed Jaguar, the 25 TF Cray XT3. Jaguar uses the AMD Opteron™ processor and has 5,212 compute processors linked in a three dimensional torus architecture using the Cray developed SeaStar™ interconnect. SeaStar provides a very high bandwidth link between processors with low latency, exactly what is needed for scalability of applications. Jaguar will be upgraded to a 250 TF system by the end of 2007 with two intermediate steps in 2006. The first upgrade will double the total system performance to 50 TF by replacing the Opteron™ processors on each of the compute nodes with AMD's latest dual-core Opteron™ processor. The dual-core processor will logically appear to be two separate CPUs. This step allows the users to begin scaling applications to over 10,000 threads of execution. In late-2006, an additional 68 cabinets will be added to the existing 56 cabinets. This will further double the system size to 100 TF. These new 68 cabinets will consist of Cray's "Hood" compute blades. Hood is the internal code name for the newest Cray

board.  It uses a newer version of the dual-core Opteron™ that supports faster DDR2 memory (62% more bandwidth than DDR1), and uses Cray's second generation SeaStar™ interconnect chip.  The SeaStar2 doubles the bandwidth between the Opteron™ and the SeaStar™, while maintaining compatibility with the SeaStar1 network.   The resulting system will have 23,016 CPU cores and give users a system for breakthrough science while continuing to scale applications to 20,000+ threads of execution.  The final upgrade for Jaguar will come in late 2007.  At this time, Cray will replace the dual-core processors in the 68 "Hood" cabinets with a future AMD processor that will take the machine to over 250 TF in the compute partition.  The initial 56 cabinets will remain as dual-core processors and will be used as a data analysis partition.   Codes running in this partition will have full access to the high speed network and file system for pre- and post-processing of data.

In 2008, the LCF will install a 1 PF computer system from Cray that is code named "Baker."   Baker continues the architecture of the XT3 and Hood systems to ease the job of porting applications.  The Baker system will continue to use a future AMD Opteron™ processor, but will use a faster memory technology and will increase the bandwidth between the processor and the interconnect by using a new version of the HyperTransport technology that AMD uses to link processors with external devices. The result will be a system that maintains the current ratios of memory, processor, interconnect, and I/O, while greatly increasing the total compute capability for the users of leadership computing.

Beyond 2008, the LCF will continue to provide leadership systems for the nation. While it is always difficult to predict the future, predicting technology trends five years out is especially perilous.   The DARPA High Productivity Computing Systems (HPCS) project gives us an unusual glimpse into what may be the first system to achieve a sustained petaflops performance on an application.  The Cray roadmap for HPCS shows a hybrid system that includes Opterons, vector processors, multi-threaded processors, and reconfigurable processors, all sharing a common infrastructure for job launch, I/O, and other services.  This is one vision of a follow-on system for delivery in 2010 or 2011.

## 5.  Opportunities and Challenges

The complexity of going from 25 TF to 1 PF in three years cannot be overstated.   The high-performance computing industry and the science community have been focusing on scaling systems to ever larger sizes for many years.    While the challenges are daunting, careful analysis and planning allows us to break up the problems into manageable sections that can be individually overcome.  The LCF has broken the tasks down into four major areas.

**Application Scaling** is the hardest of the tasks.  The SciDAC program has provided an exceptional start, but there is much work to be done.  The renewal of the SciDAC program, known as SciDAC2, will provide a solid base of funding and people to continue to scale applications and problems up to the petaflops systems we will deliver.  However, this is just a part of the needed work.  As described above, we have a series of increasingly larger machines to provide platforms for the scaling work.   A significant portion of these systems must be reserved and used for developing and testing the scaling of applications.  Further, the LCF will provide tools and test platforms to begin tuning applications to use the dual-core, and future Opteron™ processors.  The LCF will also begin assessing the scalability of critical libraries and solvers and work with the developers of these tools to scale them to the sizes of the leadership systems.

**System Software** is always one of the last items finished in the development of a new computer system.   Of course, the system

software cannot be completed and fully tested until the final hardware is delivered, but the LCF is working with Cray and others to insure that the operating system, file system, programming environment, libraries, and other system software are delivered in working condition with the delivery of the hardware. A key early milestone for Cray is to develop a plan for the internal and external resources that they will need to develop and test the system software. In most projects of this scale, the hardware is delivered late and the software development team is not given the required time to complete, test, and mature the software for the system. The LCF and Cray are putting a very high priority on both meeting the hardware delivery milestones, and in providing adequate simulators and other test platforms so that the system software development team has the resources needed to deliver their work on time.

**Space, Power, and Cooling** are rapidly becoming a major bottleneck to delivering major computer systems. Today, the LCF is using 2 megawatts (MW) of power and over 600 tons of cooling. The Baker system in 2008 will increase the total power consumption to almost 12 MW with the need for 3,600 tons of cooling. Projections for power requirements for HPCS systems range from 40-60 MW! There are very few sites in the country that can provide this much power and cooling, as well as the space needed to house such a system.

The LCF is housed in a 40,000 square foot computer center located on the ORNL campus in Oak Ridge, Tennessee. The facility today has over 11 MW of power available, with construction beginning on further upgrades to increase the power to 40 MW by 2010. The laboratory will complete construction in January 2007 of a new 70 MW power substation that has expansion capability to go to 170 MW by adding additional transformers. The computer center today has 3,600 tons of cooling, and this will be cross connected to two other

chiller plants on campus, providing as much as 15,000 tons of cooling and an upgrade path to bring in more capacity as needed.

**Infrastructure systems** provide the environment in which the leadership systems function. Scaling these systems for the anticipated load from a 40x increase in computing power requires careful planning and a solid, scalable architecture. The LCF has an architecture based on highly scalable technologies. The archival storage system, HPSS, has already demonstrated scaling at several sites to the required data volume and a substantial fraction of the required bandwidth. As a development partner for HPSS, we will focus our efforts on testing and validating that HPSS will be ready for the demand and on adding the required hardware. The LCF external networking infrastructure is also highly scalable, with the capability to increase bandwidth as needed. The internal networking was just upgraded to 10 Gigabit Ethernet. This technology is scalable to provide bandwidth as needed. The least scalable technology in use in the LCF today is the shared file system that links our computers with visualization and data analysis capability. Today, we use NFS, which is both slow, and not scalable to large numbers of systems. The LCF has initiated a project to develop and deploy a scalable shared file system for the center based on the Lustre file system from Cluster File Systems, Inc. The first version of this system will be deployed in 2006 with further expansions coming in 2007 and 2008. Our data analysis systems are being built around cluster technology. Today, the LCF has two clusters available to users for this work. We will upgrade these clusters as demand grows. In 2007, we will dedicate a 50 TF portion of Jaguar for data analysis to test the utility of a fully integrated analysis partition as an alternative to external clusters.

## 6. Summary

The LCF is well positioned to field leadership systems for the open science community in the United States. We have a strong hardware roadmap, a plan to deliver the required system software, a DOE program and LCF work to deliver the applications, and a world class facility to house the project. The challenges are large, but we are uniquely positioned to succeed in delivering breakthrough science for the nation.

## References

1. *Top 500 Supercomputer Sites*, Dongarra, Meuer, Strohmaier, Simon; A copy of the latest report may be downloaded from: http://www.top500.org
2. *Federal Plan for High-End Computing, Report of the High-End Computing Revitalization Task Force*, (National Coordination Office for Information Technology Research and Development, May 10, 2004). This document can be downloaded from: http://www.nitrd.gov/pubs/2004_hecrtf/20040702_hecrtf.pdf
3. *A Science-Based Case for Large-Scale Simulation*, (Office of Science, U.S. Department of Energy), Volume 1 (published, July 30, 2003), Volume 2 (September 19, 2004). These volumes may be downloaded from: http://www.pnl.gov/scales/
4. *Scientific Discovery through Advanced Computing*, (Office of Science, U.S. Department of Energy, March 24, 2000). A copy of this report may be downloaded from: http://www.scidac.org/SciDAC.pdf
5. *The Challenge and Promise of Scientific Computing* (Office of Science, U.S. Department of Energy, December 2003). http://www.science.doe.gov/Sub/Occasional_Papers/1-Occ-Scientific-Computation.PDF
6. *American Competitiveness Initiative – Leading the World in Innovation,* (Office of Science and Technology Policy, February 2006) A copy of this report may be downloaded from: http://www.whitehouse.gov/stateoftheunion/2006/aci/aci06-booklet.pdf
7. *High Performance Computing and Communications Act of 1991*, Bill S.272; A copy of this bill may be viewed at: http://thomas.loc.gov/cgi-bin/query/D?c102:4:./temp/~c102yqM4cq::

## Acknowledgements

## About the Authors

Arthur Bland is the Director of Operations for the NCCS and a 26-year veteran of the High-Performance Computing business. He can be reached by email at BlandAS@ornl.gov. Ann Baker leads the HPC Operations Group of the NCCS and can be reached by email at bakerae@ornl.gov. Shane Canon leads the Technology Integration group of the NCCS and can be reached by email at canonrs@ornl.gov. Al Geist is the acting Chief Technology Officer of the NCCS and leads the Network and Cluster Computing group in the Computer Science and Mathematics division at ORNL. Ricky Kendall leads the Scientific Computing group of the NCCS and can be reached by email at kendallra@ornl.gov. Jeff Nichols is the Acting Director of the NCCS and the

Division Director of the Computer Science and Mathematics Division at ORNL. He can be reached by email at nicholsja@ornl.gov. Julia White leads the User Assistance and Outreach group of the NCCS and can be reached by email at whitejc@ornl.gov. Thomas Zacharia is the Associate Laboratory Director for Computing and Computational Sciences at ORNL and may be reached by email at zachariat@ornl.gov.