

# LACSI Impact on ASC Projects at LANL

February 7, 2006

Most of the LACSI projects are driven by long-term goals. However, these projects are structured to deliver intermediate results that have impact in both the short and medium term. Many of these intermediate payouts have had significant direct impact on the ASC weapons programs; others will bear fruit in the near future. It is a goal of this report to elaborate these impacts.

The material in this report is structured according to a series of projects, with the impact of each project detailed separately. In some cases, particularly with the more mature projects, these impacts have already been realized. In other cases, the research projects have been redirected to deliver future impacts, based on our collaborations with LANL staff members.

## Performance Analysis Tools

Rice computer scientists have designed and implemented HPCToolkit, a suite of novel performance analysis tools that are dramatically simpler to use and provide a superior user interface for analyzing multiple factors affecting the node performance of scientific programs.

These tools support effective analysis of the node performance of large-scale scientific applications consisting of thousands of procedures, hundreds of thousands of lines of code, and external (possibly binary-only) libraries. The toolkit is designed to work directly with optimized application binaries. By doing so, the toolkit is language independent and it avoids the need for manual instrumentation, changes to the build process, and recompilation. A feature that distinguishes HPCToolkit from any other performance analysis tool is that it analyzes application binaries, recovers information about the loop nesting structure in the program, and maps performance back to the source code at the loop level—the level that matters for optimized scientific programs. The toolkit is multiplatform (including Alpha+Tru64, MIPS+IRIX, Pentium+Linux, Opteron+Linux, and Itanium+Linux)—a range of architectures of past, current, and emerging interest to ASC in general and LANL in particular.

HPCToolkit supports scalable data collection using hardware performance monitors for analysis of both serial and parallel codes. To facilitate interpretation of performance data, toolkit components combine measured metrics with program structure information from the application binary and correlate it with the application source code. The toolkit includes an intuitive browser for interactive exploration of application performance data. A novel and important feature of the performance browser is its support for top-down analysis, which is essential for effective analysis of large codes.

Recent work at UNM is designed to complement the performance monitoring work being done at Rice University. While Rice's HPCToolkit focuses primarily on CPU performance on a single machine, new work begun at UNM is focused on fine-grained, lightweight measurement of operating system and networking software performance. Recent work at CCS-3 has shown that operating system performance is critical to exposing the full capabilities of large-scale machines to ASC applications such as SAGE; Prof. Patrick Bridges at UNM has recently begun work on message centric monitoring, a new performance monitoring system designed to measure OS and message passing performance at runtime with minimal overhead. This system is designed to enable the deployment of operating system and runtime components that can adapt their behavior to the varying needs of ASC applications. In the past year, the UNM group demonstrated that message centric monitoring introduces minimal overhead and can be used to detect static load imbalance.

As ASC systems continue to grow toward petaflops, system power consumption continues to rise; LANL's new facility has a power budget of several megawatts. This power budget is both a critical design constraint and an operational limitation for LANL facilities and software resource management. With the increasing scale of leadership-class computing systems, node failure has become inevitable. Detecting and coping with failure is becoming critical for system reliability. To characterize power consumption and assess the general health of nodes in such systems, the UNC team has developed a system monitoring library, Health Application Programming Interface (HAPI), with a standard interface for acquiring temperature readings from processors and motherboards, voltages from power supplies, disk drive error rates, fan speeds, and other metrics. HAPI supports a wide range of data collection interfaces, allowing integration of data from diverse sources. The team has integrated the HAPI library into the SvPablo performance analysis toolkit, aiming to enable temperature and power consumption profiles of application execution and to help guide ASC application tuning, reduce LANL system power demands, and increase system reliability.

To develop efficient software environments, understanding the dynamic behavior of parallel programs is essential. Event tracing is the standard approach for obtaining the detailed data that describes application execution dynamics, however, it generates an unacceptable amount of data for ASC-class applications. To retain the benefits of event tracing with lower overhead for long running and massively parallel ASC applications, the UNC team has been integrating a compact application signature modeling method with the SvPablo toolkit. The goal is to generate performance signatures that summarize time varying application behavior during the execution of an instrumented ASC application and provide a graphical user interface for signature display and comparisons.

**Training.** In 2003, we held two performance analysis workshops at LANL to train LANL personnel on how to use LACSI-supported tools for performance analysis of their codes. LACSI researchers provided hands-on training to members of the Telluride, Shavano, Crestone, and Eolus code teams, as well as members of the transport group from X division. In the first workshop, researchers from Rice, UH, Illinois, and

Tennessee worked with LANL scientists on codes including Truchas, SAGE, and FLAG. The Eolus team was represented in the sessions, but declined to study their code at the workshop because of export control concerns. In addition to the workshop attendees, the Rice team also assisted the Partisan team in using HPCToolkit to study their code's performance in an extended teleconference.

In the second performance analysis workshop, teams from Rice (HPCToolkit) and UNC/Illinois (SvPablo and DynaProf) returned to LANL to work with broader groups from the Crestone and Shavano teams. Each tool team was teamed with a code team for a full day and we switched for the second day. SvPablo complements HPCToolkit by enabling source code instrumentation of application codes and correlation of hardware and software metrics, together with support for dynamic application adaptation. Based on experiences in the workshops and on our interactions with LANL application researchers, the Illinois/UNC team enhanced SvPablo to simplify instrumentation of multi-file codes with integrated build procedures.

The Tennessee team worked with Harvey Wasserman, Olaf Lubeck, and Michael Lang on the PAPI installation on the Opteron system as well as answered questions about performance monitoring hardware.

**Scalability.** The combination of stockpile stewardship requirements, which require high performance, and the emergence of multicore chips, each with multiple processors, further exacerbates the challenge of scalable performance monitoring and application optimization. Parallelism levels of several thousand nodes are now possible and tens to hundreds of thousands are near. Traditional task-level performance monitoring can be invasive and adversely affect the performance of the applications being measured. Moreover, monitoring each processor can produce prodigious amounts of performance data. To retain the benefits of detailed measurement while reducing data volumes, the UNC team has developed the Adaptive Monitoring and Profiling Library (AMPL), a general-purpose toolkit that reduces collected data using clustered population sampling techniques. AMPL also provides facilities for stratification of clusters into groups of monitored nodes, which can further reduce monitoring overhead if the variance of monitored metrics in these groups is low.

Experiments with the sPPM code, when collecting common PAPI hardware performance counter data on 256 node systems, show that only 5% to 14% of the total number of nodes need to be monitored to guarantee 90% confidence and 8% error in measuring the mean value. UNC also used AMPL to demonstrate the potential of low-variance stratification schemes for further reduction of monitoring overhead. For 99% confidence and at most 1% error of measurement, we have found that stratification can further reduce monitoring overhead for L2 cache misses by as much as 70% over unstratified schemes.

**Deployment for Production.** The Rice team has worked with members of CCN-8 (including Brett Kettering, Susan Post, Chip Kent, and David Montoya) to deploy HPCToolkit on all major platforms at LANL. To date, it has been deployed on open and classified SGI Origin, AlphaServer SC, and Opteron systems. The Rice team is working with CCN-8 to enhance the capabilities of HPCToolkit to provide superior support for performance diagnosis on Clustermatic Opteron clusters. In particular, the Rice team is integrating HPCToolkit with kernel support for system-wide performance data collection.

HPCToolkit is being used by LANL for performance evaluation of stockpile stewardship applications. We were informed that it was used for analysis of FLAG for the ASC burn code review in August 2003. (See the “Analysis and Tuning of ASC Applications” section.)

In the summer of 2005, the Rice team worked with David Montoya to diagnose problems that LANL was encountering when deploying HPCToolkit on Clustermatic systems. Through a joint investigation, the problem was pinpointed as an OS issue rather than an HPCToolkit issue. LANL was to follow-up with the OS deployment team.

In January 2006, Mellor-Crummey built a new version of HPCToolkit on QSC and provided sources and binaries so that the newest version of the tools could be deployed on classified systems at LANL for use by the X-8 performance analysis team.

## **Adaptability, Scalability and Fault-Tolerance**

Work at Illinois (and now at UNC) and University of Tennessee has focused not just on development of traditional performance tools, but also on the next generation of performance scalability and fault tolerance problems to be encountered by LANL and ASC applications. With emerging systems containing thousands (and soon tens of thousands) of processors, new approaches are needed that can estimate performance efficiently, enable adaptability in the face of component failures, and ensure high application performance.

**Reliability and Fault Tolerance.** Operational experience with LANL’s systems has shown that they are sensitive to transient faults, due to radiation (e.g., single bit upsets) and software glitches. Simply put, the standard assumption that system hardware and software are fully reliable is much less credible, particularly in LANL’s high altitude environment with very large systems. Concomitantly, understanding application sensitivity to system failures is critical to establishing confidence in the outputs of stockpile stewardship applications and in specifying systems and architectures that can continue operation given component failures.

The UNC team is developing reliability tools such as the Health Application Programming Interface (HAPI) for fault monitoring and the Failure Indicator Toolkit (FIT). FIT provides advanced failure detectors that enable adaptation decisions, based on a suite of libraries for analyzing computer system health, and can detect the likely failure of the monitored systems. The team is also exploring methods of building failure models

based on collected system health data to make online failure predictions. In early experiments, the failure detectors have been tested and evaluated on failure data for hard drive errors from external sources. Using ASC benchmarks, the experiments have demonstrated that relatively simple techniques can be used to monitor and predict system failures. The ability to forecast failures will help make scheduling and adaptation decisions based on reliability in addition to performance.

The Tennessee team has been working extensively with Rich Graham and others at LANL to develop an OPEN-MPI implementation.

**Fault-Tolerant Algorithms.** The Tennessee team has been working on fault-tolerant methods for iterative solutions of large sparse systems of equations. The goal of the FT-LA project is to discover and evaluate algorithms and techniques for fault-tolerant linear algebra operations on massively parallel computing environments, and encode the results in software tools that facilitate the research and help disseminate its benefits. Accordingly, FT-LA's research program divides roughly into three main parts. The FT-LA project aims to

1. *Discover and develop algorithmic innovations*, focusing on novel "lossy" algorithms that avoid checkpointing and its costs altogether but still reach the solution with reasonable performance, and on coding algorithms for algorithm-based diskless checkpointing scenarios that can succeed in massively parallel situations, where they have previously been unusable.
2. *Analyze parameter-driven performance and error propagation factors* of our new algorithms, enabling accuracy and performance through prediction and parameter selection, and leading to simulation models that support the automatic choice of optimum runtime parameters.
3. *Develop and disseminate software for creating fault tolerant numerical libraries*, facilitating both FT-LA research and enabling the community to utilize its positive results to implement resilient versions of key modules and to integrate them with libraries of choice.

**Adaptability.** Given the sensitivity of LANL systems to transient errors, the UNC team has been exploring adaptation tools and techniques for LANL and ASC applications. The goal is to develop tools and methodologies that can help applications balance high performance and reliability requirements in the face of variability and failures in the application and system components.

With additional funding from another project (NSF VGrADS), UNC has been developing a multilevel fault tolerance API. The API allows users to clearly articulate performance and reliability expectations in high-level terms. These high-level application expectations are then translated to specifications that can be realized within the resource space. At the lower level, tools and techniques such as HAPI are used to obtain health-related diagnostic information. The adaptation can be driven by fault detectors and predictors from FIT. The adaptation at the cluster resource layer needs to be also coordinated with any higher-level adaptation at the workflow or application level.

Given the multi-megawatt power consumption of LANL and other ASC systems, we have also been investigating job scheduler power optimizations. The goal is to maintain a system at a desired power level that can be specified as a percentage of power usage between the idle power and the system's maximum power. We have compared and evaluated traditional scheduling policies, including the standard First Come First Serve policies with and without backfill (FCFS and FCFS-BF), with our policies that give low power jobs higher priority over high power jobs with and without backfill (POWER and POWER-BF).

We have observed that for different power thresholds, the POWER and POWER-BF policies give a throughput lower than FCFS does, but higher than FCFS-BF gives. For low power jobs, when power is maintained at the desired level, the average and maximum delay of the jobs that are below the power threshold is considerably lower than that of the FCFS scheme. This suggests implementing specific power policies for adaptation in some control centers.

## **Analysis and Tuning of ASC Applications**

Rice computer scientists have used ASC applications to drive research in compiler technology and run-time libraries for large-scale scientific codes. This work has two goals: (1) to identify long-term research challenges in algorithms, data structures, and compiler technology that are motivated by problems faced by ASC applications and other large-scale scientific codes, and (2) to identify important code improvements that can be performed manually in the near term to improve performance of the ASC workload. In the course of this work, we studied several unclassified ASC applications, identified opportunities for improving application performance, and reported our results back to code teams at LANL. Several of these improvements have been particularly noteworthy. We mention a few of our direct impacts on ASC applications below.

**SAGE.** Rice computer scientists identified a subtle but serious performance degradation in SAGE that resulted from unnecessary copying of on-node data. The copying was a result of the strategy employed by the code for managing off-processor data. We performed experiments that demonstrated substantial performance improvement by removing this bottleneck. We advised members of SAGE team (Mike Gittings and Tom Betlach) about a strategy for restructuring SAGE's handling of off-processor data and its benefits. Improvements that they implemented roughly doubled solver performance on Blue Mountain, and boosted it by roughly 50% on ASCI Q.

A further study of SAGE by the Rice team also identified that sparse-matrix-vector product computation, a key component of many ASC codes, tends to be inefficient on modern microprocessor-based systems. Working with a benchmark derived from SAGE timing tests, Rice computer scientists identified an interaction between sparse matrix representation, sparse matrix data, and compilers that was responsible for low sparse-matrix-vector product performance in SAGE. We devised a new sparse-matrix representation that makes it possible to reorganize sparse matrix computation for higher efficiency. Working with Harvey Wasserman at LANL, we demonstrated that this new

sparse-matrix-vector multiplication improves sparse-matrix vector product performance in SAGE (with on-node data) by a factor of 2 on Itanium2 and by 47-71% on Power3-II. Computation on off-node data could benefit by applying the same technique.

**Sweep3D.** As a transport application, sweep3D represents an important part of the ASC workload. At Rice, it has served as a driving application motivating the development of performance tools including a memory hierarchy simulation package and the HPCToolkit performance tools. Although the memory hierarchy simulator was principally built for internal use for studying applications to gain insight into opportunities for compiler-based improvement, it has also seen external use at Sandia. Sweep3D has been the subject of several investigations into program transformations for improving memory hierarchy performance. As a result of this work, several higher performance variants of Sweep3D were provided back to LANL with a report showing 20% single-processor speedup. Further investigation uncovered additional opportunities for applying program reordering transformations (in particular, loop fusion). The transformed version of this application delivers 44% higher single-processor performance on Alpha systems on a  $150^3$  problem size. The same version delivers 90% improvement on single processors of the SGI Origin.

**Blanca.** The Blanca project was encountering a severe performance bottleneck in AMR vertex and zone setup. Jay Mosso and Richard Barrett produced a representative standalone unclassified code that showcased the problem. Analysis by the Rice team showed that the algorithm and data structure for maintaining an ordered collection of vertices was a severe bottleneck for large problem sizes. We implemented two alternative strategies – one for maintaining an ordered collection on the fly and a second for sorting and eliminating duplicates. The algorithm and data structure changes improved the performance asymptotically from  $O(n^2)$  to  $O(n \log n)$ , which had the effect of improving performance by a factor of 26 for the 5-level refinement test case.

**CHAD.** CHAD uses an irregular mesh data structure at the heart of its computation. Rice computer scientists analyzed the impact of data ordering on cache performance in CHAD. We showed that intelligent data orderings based on space-filling curves can double performance over naïve orderings for irregular and adaptive applications. Experiments with the volume module in CHAD showed that inadequate loop fusion and ineffective scalarization of Fortran 90 array notation increased the memory traffic by 40%. This experience was the motivation for compiler research by Rice computer scientists into better algorithms for scalarization of Fortran 90 array operations, loop fusion, and array contraction to reduce memory traffic.

**Classified applications.** In November 2005, Ken Kennedy worked at LANL with the X-8 group to analyze classified ASC codes using HPCToolkit. Kennedy and the X-8 team used the Pentium+Linux and Alpha+Tru64 versions of HPCToolkit to study performance issues of classified ASC codes on different computing platforms. Appropriately cleared personnel can obtain details about this effort and its findings from either Kennedy or Hank Alme in a classified setting. One unclassified finding of this effort was that it would be incredibly helpful to associate run-time costs, including communication costs, with the

full dynamic context (i.e., call stack) in which they are incurred. Rice had anticipated this need and has been performing research into techniques that can enable multiplatform tools to provide this capability efficiently for unmodified optimized code. In 2005, Rice developed a prototype of a portable call-stack profiler and experimented with it on both Alpha and Opteron platforms. At present, the prototype call-stack profiler requires some compiler-recorded information to enable it to unwind call stacks at *any* point in an execution. Unfortunately, most or all compilers don't record sufficient information to enable call-stack profiling of all optimized codes. Rice is presently exploring strategies for recovering the necessary information directly from application binaries to enable this precise and efficient profiler to be used on arbitrary optimized executables generated by any compiler.

At a meeting in the fall of 2005, the LANL X-8 group described having difficulties diagnosing storage management issues for large-scale ASC codes. Rice is presently working to deliver a modified version of their call-stack profiler that quantifies the amount of memory allocated and deallocated in each calling context. This tool will make it possible to pinpoint elusive memory leaks in production-scale executions.

## **Performance Modeling**

The Performance and Architecture Laboratory at LANL has developed powerful techniques for modeling the scalability of parallel codes on large-scale systems. One limitation of their approach is that they lack the ability to accurately predict node performance of whole complex codes on systems that are unavailable – this includes systems such as the Earth Simulator as well as future systems such as processor-in-memory systems. Rice computer scientists have developed a complementary strategy for semi-automatically synthesizing cross-architecture performance predictions for whole programs. Prediction tools developed at Rice have been used to make accurate cross-architecture performance predictions (e.g. accurately predict performance on MIPS and Itanium systems given only SPARC binaries) and to study the impact on performance of changing memory hierarchy characteristics.

A key feature of our strategy for developing models for predicting node performance is that the predictions are scalable. Experiments with modest problem sizes can be used to accurately predict node performance characteristics for significantly larger problem sizes that may not be feasible to simulate (or even to run) on today's systems. A Rice student interned with the PAL group during the summer of 2003. Our aim is to integrate the complementary technologies developed at Rice and LANL to better support scalable cross-architecture predictions. A key goal of this work is to influence the design of emerging petascale machines by projecting how well their proposed architectures will support ASC workloads.

Currently the Rice performance modeling and analysis team is interacting with Olaf Lubeck in the refinement of measurement and analysis for quantifying the impact of memory latency on the performance of ASC codes executing on modern microprocessor-based systems.



## **Open-Source Compiler Technology for ASC Applications**

Economic forces in the marketplace mean that relatively little attention is paid to the performance of scientific applications on commodity microprocessors; the emphasis of ongoing development is targeted primarily at supporting business applications. As a result, there is little focus on research and development of compiler technology to boost the performance of scientific codes. In part, this also hinders technology transfer because vendors have little interest in incorporating new algorithmic techniques in their compilers unless these techniques will boost performance of commercial workloads in addition to scientific ones. For this reason, it is necessary for universities to directly transfer compiler technology for improving the performance scientific codes to end users – government laboratories in particular – through open source compilers.

There has been considerable interest from the DOE in general and LANL in particular in supporting development of open source compilers. The NNSA OSSODA RFI is a recent byproduct of that interest. Rice was involved in a response to that RFI and has been active in subsequent meetings with the DOE steering committee to formulate a plan that will best meet the needs of the NNSA and make it possible to deploy results of computer science compiler research for use on production codes at the national laboratories.

Building source-to-source program translators is an effective way to transfer compiler technology that is developed in the course of university research on high-level high-impact optimizations for scientific programs. For this reason, Rice computer scientists have been working to create an open-source software platform that can transform production Fortran 90, C++ and C code based on the Open64 compiler infrastructure. Rice University is leading a collaborative effort that now has active developers at Rice Argonne National Laboratory, and Lawrence Livermore National Laboratory, in addition to international collaborators at the University of Vienna and the Universitat Politècnica de Catalunya. As this system has matured, Craig Rasmussen (LANL ACL) has expressed interest in using it for program transformation. Also, the Fortran 90 source-to-source transformation infrastructure we have developed in this project is supporting development of Fortran 90 support for automatic differentiation (described in the next section).

With supplementary funding from the DOE Office of Science, the Open64 compiler infrastructure is being used as the infrastructure supporting research and development of emerging programming languages including Co-array Fortran, which is under development at Rice, and Unified Parallel C (UPC), which is being developed by our collaborators at Berkeley. Ultimately, petascale machines will need to be programmed differently than today's parallel systems. Our ongoing research efforts to build compilers for global address space languages are intended as steps towards having compilers manage data movement and synchronization implicitly and simplify parallel programming for application scientists. There has been interest in Co-array Fortran at LANL, in particular from the Crestone code team. Experiments by Cray using SAGE on the Cray X1 required converting parts of it to Co-array Fortran for best performance. Experiments carried out at Rice in 2004-2005 with a Co-array Fortran implementation of Sweep3D show that Co-array Fortran can deliver MPI-level performance, yet enable

simpler global address space based programming. These experiences with Sweep3D are driving ongoing exploration of Co-array Fortran language support to simplify high-performance communication in pipelined computations.

## **Automatic Differentiation**

LACSI has supported ongoing development of the ADIFOR automatic differentiation tool. This tool has been employed by a number of projects at LANL. Primary, it has been used for verification/validation and parameter identification/estimation. Rudy Henninger's group (currently CCS-2) has used ADIFOR on MESA 1D & MESA 2D (Armor/Antiarmor codes), "Caramana's Lagrangian Test Code" (used to study hydro methods that are implemented in the Shavano project's FLAG code), and Truchas-1D (Telluride project). Ralph Nelson (Shavano Project) has used ADIFOR on a light-water-reactor safety code called TRAC. In addition to Henninger and Nelson, there are six additional LANL scientists who have registered for copies of the ADIFOR code. Three of them indicated that they obtained a copy of ADIFOR after seeing a talk at LANL by a Rice computer scientist. Current work is focused on applying the latest ADIFOR to Truchas 3D and to FLAG. A prototype forward mode version of the Adifor90 has been installed at LANL on CCS-2 computers. Differentiation of the Truchas code is ongoing. 25 of the Truchas routines have been verified (using difference approximations). Forward mode sensitivities for the simple regression tests in the Truchas code are expected to be complete by the end of February 2006. Differentiation of the entire Truchas program will be completed by the end of May 2006.

## **Compiler Technology for High Productivity Parallel Programming**

The principal stumbling block to using parallel computers productively is that parallel programming models in wide use today place most of the burden of managing parallelism and optimizing parallel performance on application developers. We face a productivity crisis if we continue programming parallel systems at such a low level of abstraction as these parallel systems increase in scale and architectural complexity, while decreasing in whole system reliability because of the proliferation of components. Growing awareness of the looming productivity crisis has led to the creation of DARPA's High Productivity Computing Systems program. LANL's Jeff Brown has been advocating that NNSA support a high productivity computing research thrust.

Since the inception of LACSI, the Rice team has been exploring compiler technology to support high productivity parallel programming languages. High-level data parallel languages based on a global view of data offer a dramatically simpler alternative for programming parallel systems. Programming in such languages is simpler: one simply reads and writes shared variables without worrying about synchronization and data movement. An application programmer merely specifies how to partition the data and leaves the details of partitioning the computation and choreographing communication to a parallelizing compiler. Having an HPF program achieve over 10 TFLOPS on Japan's Earth Simulator has rekindled interest in high-level programming models within the US.

For high-level models for data-parallel programming to be widely used, they must satisfy four criteria:

1. they must be expressive enough to support a broad spectrum of sophisticated parallel algorithms,
2. they must deliver performance competitive with that of hand-coded parallelizations,
3. they must be ubiquitously available everywhere from emerging tightly-coupled architectures to desktop workstations, and
4. they must deliver appropriate levels of efficiency on each class of systems.

The Rice team has been investigating locality-aware compiler technology under LACSI support since the founding of the institute. The principal products of this research have been language-independent analysis and code generation techniques for managing single-threaded parallel programming languages. These techniques are applicable to a broad class of languages ranging from MATLAB to Cray's Chapel. This research has garnered two best paper awards to date and has been selected for publication in three journal special issues (among other publications). To date, this research has made significant progress in developing compiler technology that helps address the productivity problem by transforming complex programs written in a high-level single-threaded form into sophisticated scalable codes that run with hand-coded efficiency on large-scale parallel systems. Ongoing work in this area is focused on enhancing compiler capabilities to more effectively support adaptive and multi-level algorithms. In 2005, Rice completed a study to evaluate the effectiveness of compiler technology developed by the Rice team for transforming IMPACT3D, a single-threaded HPF program developed for the Earth Simulator, into one that can execute efficiently on parallel microprocessor-based clusters that have been the focus of ASC investments. The Rice team found that their HPF compiler technology was able to generate scalable code that achieved approximately 18% of peak efficiency for the IMPACT3D application on processor ensembles as large as 1024 processors. This study showed that both data-parallel programming models and data-parallel compilers can deliver performance for data-intensive codes that is competitive with hand-coded parallelizations on modern microprocessor-based systems.

## **Optimization of Object-Oriented Programming Languages**

Several of the key LANL weapons codes include substantive code written in object-oriented languages, especially C++. It is our understanding that the Marmot project is considering using Python as well as C++ for coding. Object-oriented technologies present unique challenges for achieving the highest possible node performance. With these issues in mind, the Rice team has undertaken a study of optimizations for such languages. Although this study was conducted in the context of Java, the techniques are equally applicable to other languages, including C++. The principal focus of this work has been two different kinds of optimizations. First, a technique called "object inlining" can be used to eliminate the overheads, especially memory access costs, associated with the use of arrays of object data structures. Essentially, it replaces arrays of such objects with arrays of the instance variables for the given objects. Second, through the use of global type analysis, the overheads associated with dynamic dispatch, an important technique in

object-oriented programming, can be reduced or eliminated. We have developed a global type analysis algorithm that can determine—in many cases precisely—the actual method being called at a dynamic method invocation site. In codes with a significant amount of virtual function use, this can have a huge payoff.

We have demonstrated the effectiveness of our compiler technologies on a particle-in-cell code, Parsek, written in object-oriented style in Java at LANL. With our collaborators from Los Alamos, we have published these results in the *Journal of Concurrency and Computation: Practice and Experience* 2004. Our analysis and optimization infrastructure, JaMake, optimized the object-oriented version of Parsek to achieve a performance improvement of approximately 80% over the original code. The automatically produced code is very competitive with the hand-optimized JavaParsek.

Our original post-Java direction was to apply these transformation strategies to the Marmot project by adapting Dan Quinlan's ROSE infrastructure from Livermore to include the optimizations in JaMake. However, based on discussions with code teams inside X Division, we have decided to focus most of our effort on optimizations within the Ajax Programming System with the goal of dramatically improving the performance of applications written using Ajax. To that end, we are in process of trying to secure a copy of the export-controlled application FLAG, for study at Rice. We have proposed measures to the Rice administration that would allow the execution of licenses for such codes, including physical and information security measures to avoid deemed export, and we will shortly be seeking a license for FLAG. We are also hoping to get a copy of the Ajax translator source code to analyze and modify, although we have not been able to ascertain its classification status. By applying some of the standard optimization strategies for object-oriented languages, along with new ones developed by LACSI, we believe that dramatic performance improvements can be achieved.

## **Component Integration Systems for HPC**

Through the LACSI collaboration, the Rice team has come to understand the importance of applying sound software engineering strategies to ASC codes. One such strategy is the extensive use of pre-developed, reusable software components for structuring the code. Such components can be independently tested, providing additional reliability to applications that employ them.

Currently, most component-based systems suffer significant overheads for crossing component boundaries. This tends to drive the developers of HPC codes toward heavyweight components. Although this approach has proven successful in several HPC codes at the DOE laboratories, it precludes some important structuring strategies, such as the mixing of data structure components (e.g., sparse matrices or adaptive meshes) with functional components (e.g., linear algebra) because such integration would incur high overheads for cross-component method invocation. Instead, developers often bundle the desired functionality into the data structure component (or vice versa) making it difficult to change the data structure or to use the same data structure with different functions.

To address these problems, the Rice component integration group has redirected its research effort to focus on a new component integration strategy that preprocesses collections of components with the aim of dramatically reducing component-crossing penalties through compiler optimization. Although our initial prototypes have used Matlab as the language for application and component implementation (this is because many library developers use Matlab, or another scripting language, for prototyping), the techniques are applicable to any implementation language. In fact, components written in Matlab are automatically translated to C or Fortran before being used in an application. Preliminary performance results show that application run times can be improved significantly, in some cases by integer factors.

We are now driving this research effort by working with the developers of applications within the Marmot project to understand how the technology can be effectively applied to future weapons codes. Although the effort is still in its preliminary stages, the goal is to have significant impact on future generations of the Marmot applications, increasing the ease of implementation without incurring undue performance costs. We are also looking into the application of these ideas in the Ajax programming environment.

To date, this effort has convened one workshop on components, and has had two full-day collaboration and design workshops with members of the Marmot team. From those efforts, we have developed a plan of attack that would focus on several important goals. First, we would like to understand whether performance and memory problems due to multi-material calculations can be ameliorated by precompilation technologies that select optimized methods for different combinations of materials dynamically. Second, we are looking at widely used library collections, such as Trilinos (which to our knowledge has not yet been employed in LANL weapons codes), to see if they can be effectively used in a component integration system.

We have also used telescoping languages for research on high-performance parallel implementations of scripting languages, such as Matlab, R, and Python. We are planning an effort with Dave Higdon in D1 to parallelize V&V codes written for ASC in Matlab and R. Our research on Parallel Matlab, which has attracted supplementary support from DARPA (through NSF), has achieved linear speedups on Matlab stencil codes. Matlab is important because it is rapidly becoming the standard language for all engineers.

## **Automatic Tuning of Components and Applications**

Because of the substantive efforts involved in retuning applications when they are moved from one platform to another, LACSI researchers are involved in several projects to explore strategies for automatically tuning applications. These efforts are using experimental strategies that explore various ways to structure computationally intensive loop nests to achieve near-optimal performance on new platforms. This line of work is extremely relevant to the weapons program because it can lead to significant savings in tuning effort by application developers.

One model for the automatic tuning strategies is the ATLAS system developed at the University of Tennessee, a LACSI partner. ATLAS is a system for pre-tuning the LAPACK BLAS by extensive preliminary execution of the critical loop nest with different parameters for blocking and unrolling. Experiments have shown that ATLAS achieves performance that is very close to the best hand-tuned versions of the BLAS on a variety of platforms. This work has established that automatic tuning systems can be an effective substitute for human tuning, thus achieving a significant boost in programming productivity.

While ATLAS is focused on specific loop nests, several projects at Rice, Tennessee, and Houston are pursuing the goal of automatically tuning general applications and components. Rice researchers are pursuing the strategy of “adaptive compilation” which can be used to try, in a heuristic search, different orders of compilation passes, different compiler optimization settings, or different strategies for unrolling and inlining. This effort has already conducted the most extensive experiments on automatic tuning that have ever been attempted. The work has shown that large improvements over standard compilation can be achieved by finding the compilation order, parameter settings, and inlining strategy best suited to a particular platform.

A second group at Rice has been using LoopTool, which applies program transformations such as loop blocking and unroll-and-jam based on directives manually inserted in the program text, to explore automatic selection of transformation parameters. This work has shown that direct search strategies can reduce the search and execute time needed to pre-tune a general application to reasonable levels without significant loss of performance. Both Rice projects were reported on in papers selected for presentation at the 2003 and 2004 LACSI Symposia. In addition, the LoopTool-based system has been extended to include not only cache blocking and register reuse enhancement, but also loop fusion, which is critical for Fortran 90 programs that use array syntax. The resulting tuning system has achieved significant improvements on a variety of codes, notably speeding up Sweep3D by a factor of 1.3 on an Alpha processor like the ones used in the ASC Q machines.

The Tennessee group is extending the ATLAS work to more general component tuning while the UH effort has focused on automatic tuning of components of modest complexity that are decomposable in different ways, such as typical linear algebra library components. The technique has so far been applied to FFTs. FFTs on most data set sizes can be decomposed in several different ways that not only result in different numbers of arithmetic operations but more importantly in a range of opportunities to schedule sets of operations for maximum exploitation of the memory system as well as of functional units. The approach taken in exploiting the degrees of freedom in decomposition and scheduling has two phases: one carried out at library build time, one at run-time. At build time highly optimized “kernels” or subcomponents are created and their performance measured and entered into a data base. At run time, when actual problem sizes and data allocation are known, alternative decompositions down to the kernel level are evaluated for performance and an execution plan is generated for a particular decomposition with a schedule that seeks to minimize execution time. Various strategies that trade off the

number of options evaluated against the expected performance gain from evaluating additional options are currently being explored. Though the actual decompositions depend on the function to be performed, such as the FFT, the methodology used should apply to common mathematical operations that are decomposable. In carrying out this development, we have created a toolbox, which we call CodeLab, that should enable us to apply the technique to other important functions in ASC codes with a more modest effort than the one we have devoted to this first application of the methodology. The next intended target is support for codes using multigrid techniques.

All of the automatic tuning work is directly relevant to ASC goals, because it will lead to tools that facilitate better and faster migration of existing weapons codes to new HPC platforms.

In the interest of furthering this area of research, LACSI sponsored a workshop on automatic tuning at the 2005 LACSI symposium. Both Rice groups presented, as did the group from Tennessee. Other presenters included Kathy Yelick of LBL and Berkeley, Dan Quinlan (developer of ROSE) and Qing Yi of Livermore, Clint Whaley (original developer with Jack Dongarra of ATLAS), Kamen Yotov of Cornell, and Mary Hall of USC ISI. A mailing list was established for further discussions and correspondence. We expect that this discussion group will accelerate the pace of work in this area, which is critical to achieving portable performance on ASC codes

## **Operating Systems and Networking**

Systems researchers at UNM and Rice are working with the LA-MPI (OpenMPI) team on performance issues. These performance issues will be critical for ensuring that LANL's new generation of Linux clusters will meet performance and usability goals.

The Rice team has integrated an event-driven progress engine into LA-MPI. This event-driven progress engine improves overall responsiveness and performance of the LA-MPI TCP path. Since this work is recent, it is not yet integrated back into the main source tree. We have shown that with this new event-driven progress engine, the performance of MPI over TCP using 1Gbps Ethernet is similar to 1.2Gbps Myrinet on the NAS benchmarks. With a few additional enhancements, we believe that we can beat Myrinet.

Much of the recent work at UNM, in conjunction with Ron Minnich and Rich Graham of CCS-1, has focused on protocol optimization for message-passing runtimes for large-scale machines. Patricia Crowley is researching TCP scaling issues on large-scale systems and is attempting to reduce resource usage in TCP. In the past, TCP has occasionally been used for message passing in large-scale machines, though it is more frequently used for performance-critical tasks such as communication with NFS servers to read and write initialization data, checkpoints, and results. With processors containing thousands of nodes, limiting the resource usage of stateful protocols such as TCP while at the same time increasing their performance is essential to providing the full capabilities of the large-scale machines to ASC applications. Patricia completed an analytic model to characterize the benefits of partial protocol offload and validated this model. Patricia

completed her PhD in August and is continuing her research as a postdoc at UNM. Galen Shipman, an MS student at UNM, is researching protocols for Open MPI, in particular, protocols that can be used to support Infiniband. Infiniband is becoming an important interconnect technology in high performance computing. Recent efforts in large-scale Infiniband deployments are raising scalability questions in the HPC community. Open MPI, a new open source implementation of the MPI standard targeted for production computing, provides several mechanisms to enhance Infiniband scalability. Initial comparisons with MVAPICH, the most widely used Infiniband MPI implementation, show similar performance but with much better scalability characteristics. Specifically, small message latency is improved by up to 10% in medium/large jobs and memory usage per host is reduced by as much as 300%. In addition, Open MPI provides predictable latency that is close to optimal without sacrificing bandwidth. Galen completed his Master's thesis in November of 2005 and has accepted a position at Los Alamos National Laboratory.

**Common Hardware Infrastructure.** To streamline the sharing and deployment of research and software artifacts among the LACSI academic partners and LANL, we have begun the procurement by the partners of development systems very similar to the Blue Steel system in CCS-1. These development systems are clusters of dual-processor Opteron nodes with both Gigabit Ethernet and Infiniband interconnects. Because these systems will use the Clustermatic Linux software stack, the products of the academic research will be much more easily deployed at LANL and a wider range of academic personnel will be able to work directly with the systems of interest. Furthermore, these configurations will ensure that academic partners will be addressing the unique issues of using Infiniband as a communication fabric for high performance computing.

To provide a common platform for LACSI partners to conduct research and experiments, UNC has configured and deployed a 16-node Clustermatic (dual Opteron) Linux cluster based on the LANL infrastructure. We have deployed a suite of fault monitoring infrastructure tools on this cluster, including HAPI for general system health monitoring and commercial temperature monitoring sensors.

The temperature monitoring system on the UNC LACSI cluster includes a primary monitor (a 1U rack-mount device), rear sensors, front sensors, and ones on the inside skin of the racks. The sensors are distributed along the length of the entire rack, whereas the nodes fill roughly 3/4th's of the rack. These sensors sample every 30 seconds and record data in the logs that are updated every two days.

We conducted experiments with the sPPM and Sweep3D codes, measuring spatial and temporal temperature variations. The preliminary results show that there is a non-uniform variation of the temperature along the length of the rack. The sensors that demonstrate the biggest fluctuations are located along the middle of the rack, indicating a high heat density in that area. We experimented with different data sets with Sweep3D code and observed that the greater computing requirements result in larger temperature variations.



## **New Polyhedral-Mesh Diffusion Discretizations**

Yuri Kuznetsov of UH, working in collaboration with researchers in LANL groups T-7 and X-3, has developed a new diffusion discretization scheme on arbitrary polyhedral meshes that offers several significant advantages over previous schemes. The mimetic methodology was used to formulate this discretization. The breakthrough resulted from a key idea by Kuznetsov. It was suggested that a vector inner product on a general polyhedron be formed by using a decomposition of the polyhedron into a set of tetrahedrons in conjunction with an interior-cell vector interpolation and a vector inner product expression for a single tetrahedron. The proposed diffusion discretization scheme is second-order accurate and symmetric positive-definite, including problems with material discontinuities. It happens to be the first symmetric polyhedral scheme for the 3D diffusion equations. Previously used polyhedral schemes are non-symmetric. It was shown numerically on test problems relevant to LANL applications that the equations for the new scheme cost far less to solve than the non-symmetric equations. It was also shown that the new scheme is more accurate than the old ones. The important advantage of the new scheme is that it can be used for various types of meshes, including non-matching meshes and the meshes with local refinements (AMR meshes).

LANL researchers M. Shashkov and K. Lipnikov in the T-7 group, and S. Runnels in the X-3 group have recently implemented the discretization scheme for the diffusion equations in the FLAG code. In cooperation with other members of the X-3 group, K. Lipnikov and S. Runnels developed a parallel version of the code.

Future plans call for the extension of the method for the discretization of the diffusion type equations on polyhedral meshes with mixed cells. Such discretization methods for problems with mixed materials inside polyhedral mesh cells are in great demand at LANL and other institutions. Another challenging problem is to develop, investigate, and evaluate robust and efficient parallel preconditioners for discretization methods on arbitrary distorted polyhedral meshes.