

Performance Evaluation of Adaptive Scientific Applications using TAU

Sameer Shende^{a*}, Allen D. Malony^a, Alan Morris^a, Steven Parker^b, J. Davison de St. Germain^b

^aPerformance Research Laboratory, Department of Computer and Information Science, University of Oregon, Eugene, OR 97403, USA

^bScientific Computing and Imaging Institute, University of Utah, Salt Lake City, UT 84112, USA

1. Introduction

Fueled by increasing processor speeds and high speed interconnection networks, advances in high performance computer architectures have allowed the development of increasingly complex large scale parallel systems. For computational scientists, programming these systems efficiently is a challenging task. Understanding the performance of their parallel applications is equally daunting. To observe and comprehend the performance of parallel applications that run on these systems, we need performance evaluation tools that can map the performance abstractions to the user's mental models of application execution. For instance, most parallel scientific applications are iterative in nature. In the case of CFD applications, they may also dynamically adapt to changes in the simulation model. A performance measurement and analysis system that can differentiate the phases of each iteration and characterize performance changes as the application adapts will enable developers to better relate performance to their application behavior. In this paper, we present new performance measurement techniques to meet these needs. In section 2, we describe our parallel performance system, TAU. Section 3 discusses how new TAU profiling techniques can be applied to CFD applications with iterative and adaptive characteristics. In section 4, we present a case study featuring the Uintah computational framework and explain how adaptive computational fluid dynamics simulations are observed using TAU. Finally, we conclude with a discussion of how the TAU performance system can be broadly applied to other CFD frameworks and present a few examples of its usage in this field.

2. TAU Performance System

Given the diversity of performance problems, evaluation methods, and types of events and metrics, the instrumentation and measurement mechanisms needed to support performance observation must be flexible, to give maximum opportunity for configuring performance experiments, and portable, to allow consistent cross-platform performance problem solving. The TAU performance system [1,4], is composed of instrumentation, measurement, and analysis

*This research was supported by the U.S. Department of Energy, Office of Science, under contracts DE-FG03-01ER25501 and DE-FG02-03ER25561, and University of Utah and LLNL DOE contracts B524196 and 2205056.

parts. It supports both profiling and tracing forms of measurements. TAU implements a flexible instrumentation model that permits a user to insert performance instrumentation hooks into the application at several levels of program compilation and execution. The C, C++, and Fortran languages are supported, as well as standard message passing (e.g., MPI) and multi-threading (e.g., Pthreads) libraries.

For instrumentation we recommend a dual instrumentation approach. Source code is instrumented automatically using a source-to-source translation tool, *tau_instrumentor*, that acts as a pre-processor prior to compilation. The MPI library is instrumented using TAU's wrapper interposition library that intercepts calls to the MPI calls and internally invokes the TAU timing calls before and after. TAU source instrumentor can take a selective instrumentation file that lists the name of routines or files that should be excluded or included during instrumentation. The instrumented source code is then compiled and linked with the TAU MPI wrapper interposition library to produce an executable.

TAU provides a variety of measurement options that are chosen when TAU is installed. Each configuration of TAU is represented in a set of measurement libraries and a stub makefile to be used in the user application makefile. Profiling and tracing are the two performance evaluation techniques that TAU supports. Profiling presents aggregate statistics of performance metrics for different events and tracing captures performance information in timestamped event logs for analysis. In tracing, we can observe along a global timeline when events take place in different processes. Events tracked by both profiling and tracing include entry and exit from routines, interprocess message communication events, and other user-defined atomic events. Tracing has the advantage of capturing temporal relationships between event records, but at the expense of generating large trace files. The choice to profile trades the loss of temporal information with gains in profile data efficiency.

3. CFD Application Performance Mapping

Observing the behavior of an adaptive CFD application shows us several interesting aspects of its execution. Such applications typically involve a domain decomposition of the simulation model across processors and an interaction of execution phases as the simulation proceeds in time. Each iteration may involve a repartitioning or adaption of the underlying computational structure to better address numerical or load balance properties. For example, a mesh refinement might be done at iteration boundaries and information about convergence or divergence of numerical algorithms is detailed. Also, domain specific information such as the number of cells refined at each stage gives a user valuable feedback on the progress of the computation.

Performance evaluation tools must capture and present key application specific data and correlate this information to performance metrics to provide a useful feedback to the user. Presenting performance information that relates to application specific abstractions is a challenging task. Typically, profilers present performance metrics in the form of a group of tables, one for each MPI task. Each row in a table represents a given routine. Each column specifies a metric such as the exclusive or inclusive time spent in the given routine or the number of calls executed. This information is typically presented for all invocations of the routine. While such information is useful in identifying the routines that contribute most to the overall execution time, it does not explain the performance of the routines with respect to key application phases. To address this shortcoming, we provide several profiling schemes in TAU.

3.1. Static timers

These are commonly used in most profilers where all invocations of a routine are recorded. The name and group registration takes place when the timer is created (typically the first time a routine is entered). A given timer is started and stopped at routine entry and exit points. A user defined timer can also measure the time spent in a group of statements. Timers may be nested but they may not overlap. The performance data generated can typically answer questions such as: *what is the total time spent in `MPI_Send()` across all invocations?*

3.2. Dynamic timers

To record the execution of each invocation of a routine, TAU provides dynamic timers where a unique name may be constructed for a dynamic timer for each iteration by embedding the iteration count in it. It uses the start/stop calls around the code to be examined, similar to static timers. The performance data generated can typically answer questions such as: *what is the time spent in the routine `foo()` in iterations 24, 25, and 40?*

3.3. Static phases

An application typically goes through several phases in its execution. To track the performance of the application based on phases, TAU provides static and dynamic phase profiling. A profile based on phases highlights the context in which a routine is called. An application has a default phase within which other routines and phases are invoked. A phase based profile shows the time spent in a routine when it was in a given phase. So, if a set of instrumented routines are called directly or indirectly by a phase, we'd see the time spent in each of those routines under the given phase. Since phases may be nested, a routine may belong to only one phase. When more than one phase is active for a given routine, the closest ancestor phase of a routine along its callstack is its phase for that invocation. The performance data generated can answer questions such as: *what is the total time spent in `MPI_Send()` when it was invoked in all invocations of the `IO (IO => MPI_Send())` phase?*

3.4. Dynamic phases

Dynamic phases borrow from dynamic timers and static phases to create performance data for all routines that are invoked in a given invocation of a phase. If we instrument a routine as a dynamic phase, creating a unique name for each of its invocations (by embedding the invocation count in the name), we can examine the time spent in all routines and child phases invoked directly or indirectly from the given phase. The performance data generated can typically answer questions such as: *what is the total time spent in `MPI_Send()` when it was invoked directly or indirectly in iteration 24?* Dynamic phases are useful for tracking per-iteration profiles for an adaptive computation where iterations may differ in their execution times.

3.5. Callpaths

In phase-based profiles, we see the relationship between routines and parent phases. Phase profiles [5] do not show the calling structure between different routines as is represented in a callgraph. To do so, TAU provides callpath profiling capabilities where the time spent in a routine along an edge of a callgraph is captured. Callpath profiles present the full flat profiles of routines (or nodes in the callgraph), as well as routines along a callpath. A callpath is represented syntactically as a list of routines separated by a delimiter. The maximum depth of a callpath is controlled by an environment variable.

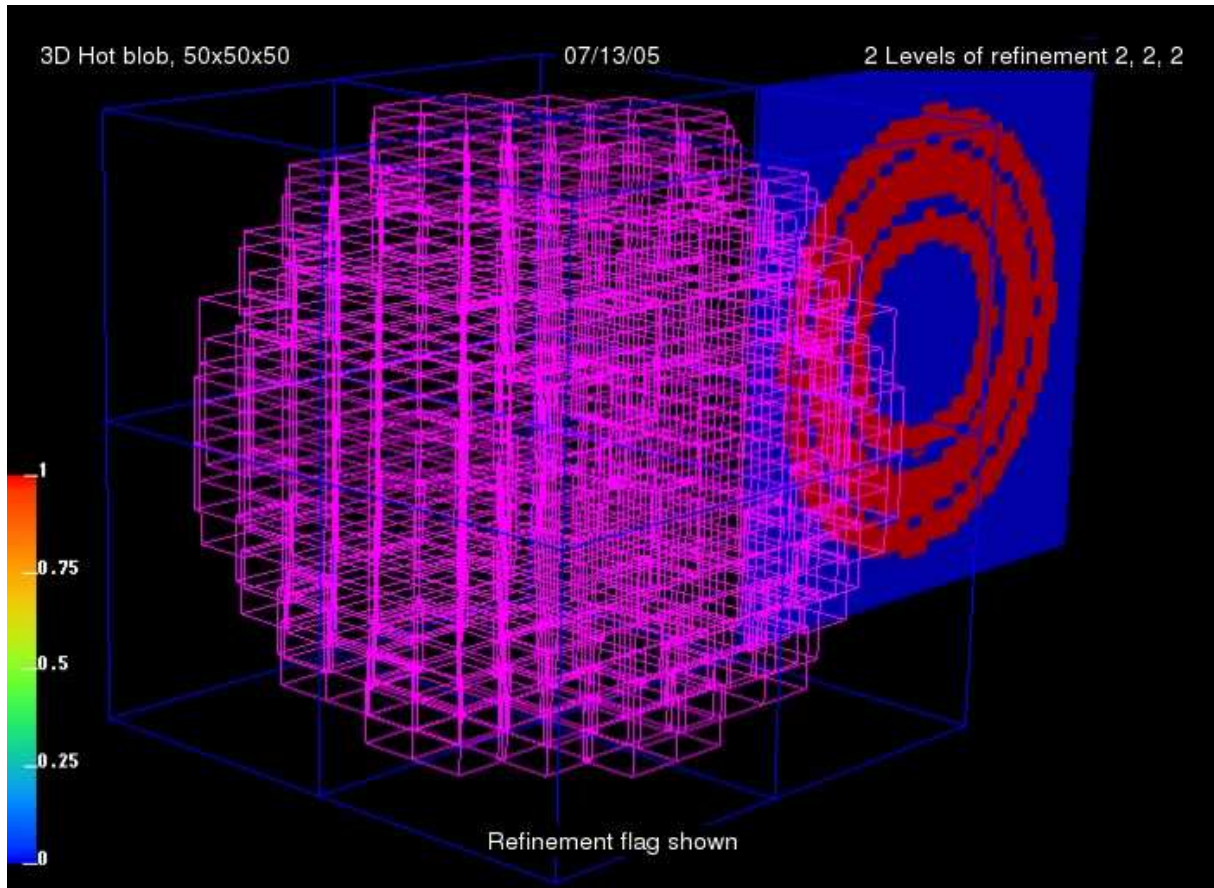


Figure 1. Adaptive Mesh Refinement in a parallel CFD simulation in the Uintah Computational Framework

3.6. User-defined Events

Besides timers and phases that measure the time spent between a pair of start and stop calls in the code, TAU also provides support for user-defined atomic events. After an event is registered with a name, it may be triggered with a value at a given point in the source code. At the application level, we can use user-defined events to track the progress of the simulation by keeping track of application specific parameters that explain program dynamics, for example, the number of iterations required for convergence of a solver at each time step, or the number of cells in each iteration of an adaptive mesh refinement application.

4. Case Study: Uintah

We have applied TAU's phase profiling capabilities to evaluate the performance of the Uintah computational framework (UCF) [2]. The TAU profiling strategy for Uintah is to observe the performance of the framework at the level of patches, the unit of spatial domain partitioning. Thus, we instrument UCF with dynamic phases where the phase name contains the AMR level

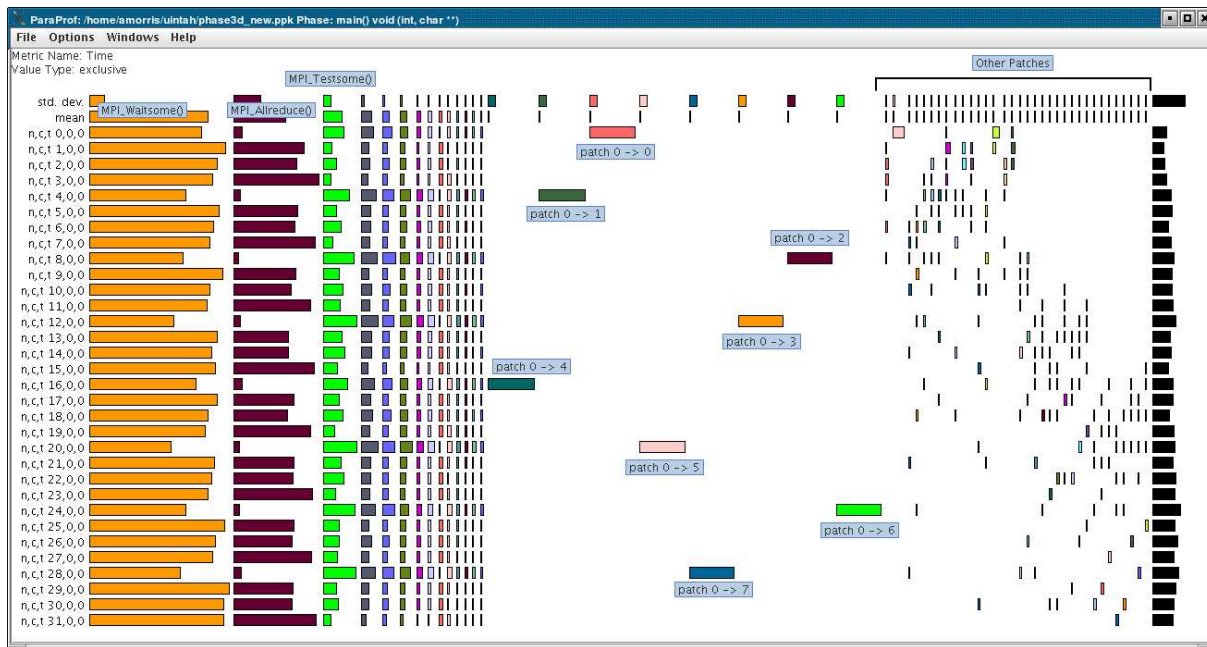


Figure 2. Distribution and time taken in various phases in the Uintah Computational Framework

and patch index. The case study focuses on a 3 dimensional validation problem for a compressible CFD code. Initially, there is a region of gas at the center of the computational domain that is at a high pressure (10atm) and temperature gas (3000K). At time = 0 the gas is allowed to expand forming a spherical shockwave. Eventually, a low pressure region will form in the center and the expanding flow will reverse directions. Figure 1 shows a sample distribution of patches in the domain. The blue outer cubes enclose the eight (2x2x2) level 0 patches. The pink inner cubes cover the “interesting” portions of the domain that have been selected by the AMR subsystem for mesh refinement on level 1.

We instrumented Uintah by creating a phase based on the level index and patch index for a task given to the Uintah scheduler. Phases are given name such as “Patch 0 ->1”, which represents the 2nd patch on level 0. All tasks and instrumented functions are then associated by way of phase profiling with the patch on which the computation is done.

Figure 2 shows the performance data obtained from the simulation as displayed by ParaProf. Here we see the eight patches that make up level 0 overwhelm the time taken in other patches. They executed on nodes 0,4,8,12,16,24 and 28. This gave us an immediate result showing that these nodes spent extra time processing the level 0 patches, while the other nodes waited in MPI_Allreduce.

Figure 3 shows the distribution of a given timer, ICE::advectAndAdvanceInTime, across all phases, each representing a patch. We can see that the time spent in this task for the eight level 0 patches is more 9 than seconds, while the time spent for all the other patches is less than 4 seconds each.

Figure 4 shows the profile for the phase ”patch 0 ->0” which runs only on node 0. We see the partition of each task that was executed under this phase. If this phase had run on other nodes,

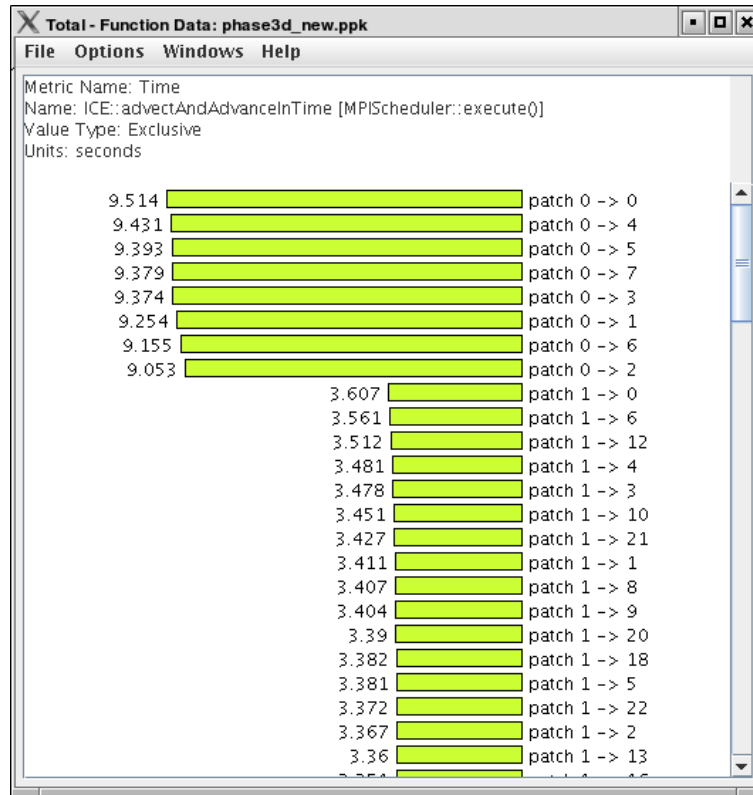


Figure 3. The distribution of the ICE::advectAndAdvanceInTime task across patches using phase profiling

we would have aggregate statistics as well (mean, std. dev.).

Phase profiling in UCF allows the developers to partition the performance data for tasks and instrumented functions across spatially defined patches with information from the AMR subsystem. This data can be used to identify load balancing issues as well as establish a better understanding of code performance in an AMR simulation.

5. Other frameworks

Besides the Uintah computational framework, TAU has been applied successfully to several frameworks that are used for computational fluid dynamics simulations. These include VTF [9] from Caltech, MFIX [7] from NETL, ESMF [10] coupled flow application from UCAR, NASA and other institutions, SAMRAI [12] from LLNL, Miranda [11] from LLNL, GrACE [13] from Rutgers University, SAGE [6] from SAIC, and Flash [8] from University of Chicago. Our work in performance evaluation of adaptive scientific computations can be broadly applied to other CFD codes. Thus, CFD frameworks can benefit from the integration of portable performance profiling and tracing support using TAU.

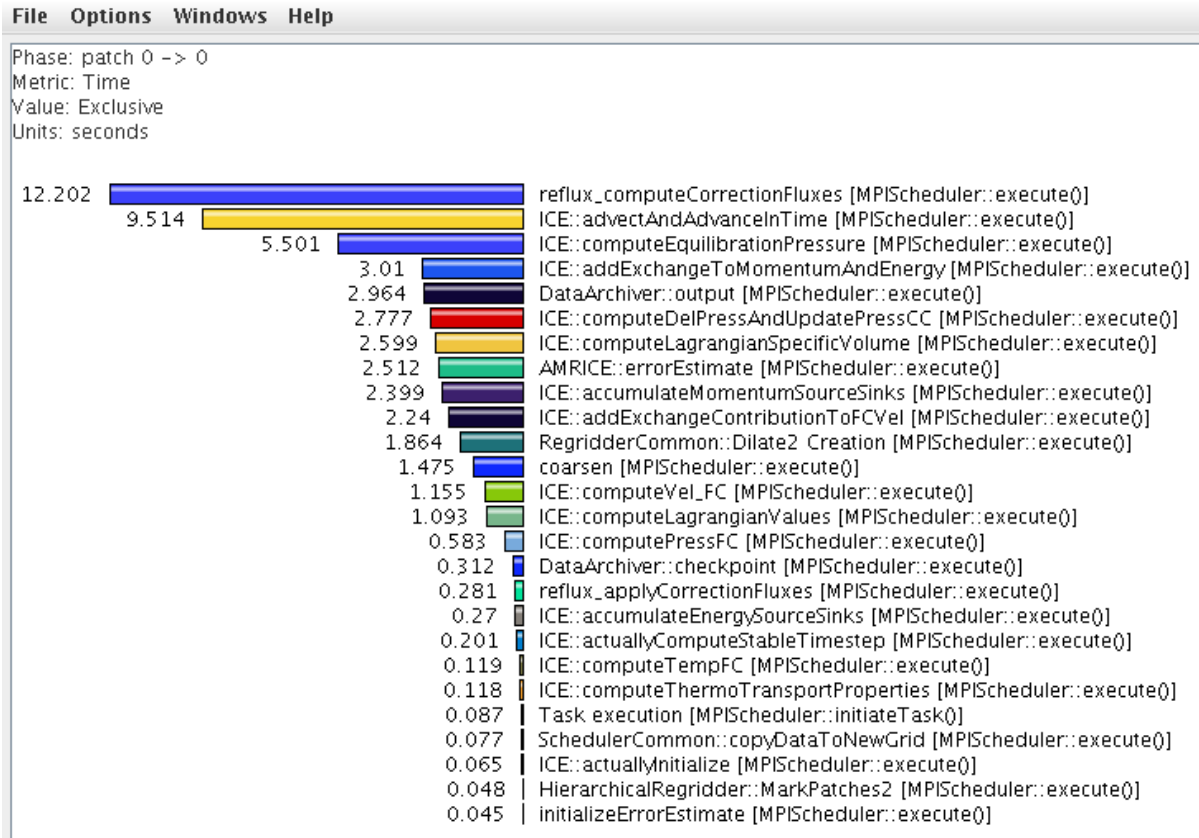


Figure 4. Phase profile for patch 0 ->0 shows tasks executing under this phase

6. Conclusions

When studying the performance of scientific applications, especially on large-scale parallel systems, there is a strong preference among developers to view performance information with respect to their “mental” model of the application, formed from the structural, logical, and numerical models used in the program. If the developer can relate performance data measured during execution to what they know about the application, more effective program optimization may be achieved. In this paper, we present portable performance evaluation techniques in the context of the TAU performance system and its application to the Uintah computational framework. We illustrate how phase based profiling may be effectively used to bridge the semantic gap in comprehending the performance of parallel scientific applications using techniques that map program performance to higher level abstractions.

REFERENCES

1. A. D. Malony and S. Shende and R. Bell and K. Li and L. Li and N. Trebon, “Advances in the TAU Performance System,” Chapter, “Performance Analysis and Grid Computing,” (Eds. V. Getov, et. al.), Kluwer, Norwell, MA, pp. 129-144, 2003.

2. J. D. de St. Germain and J. McCorquodale and S.G. Parker and C.R. Johnson, "Uintah: A Massively Parallel Problem Solving Environment," Ninth IEEE International Symposium on High Performance and Distributed Computing, IEEE, pp. 33–41. 2000.
3. R. Bell and A. D. Malony and S. Shende, "A Portable, Extensible, and Scalable Tool for Parallel Performance Profile Analysis," Proc. EUROPAR 2003 conference, LNCS 2790, Springer, pp. 17-26, 2003.
4. TAU Portable Profiling. URL: <http://www.cs.uoregon.edu/research/tau>, 2005.
5. A. D. Malony and S. S. Shende and A. Morris, "Phase-Based Parallel Performance Profiling," Proc. ParCo 2005 Conference, Parallel Computing Conferences, 2005.
6. D. Kerbyson and H. Alme and A. Hoisie and F. Petrini and H. Wasserman and M. Gittings, "Predictive Performance and Scalability Modeling of a Large-Scale Application," Proc. SC 2001 Conference, ACM/IEEE, 2001.
7. M. Syamlal and W. Rogers and T. O'Brien, "MFX Documentation: Theory Guide," Technical Note, DOE/METC-95/1013, 1993.
8. R. Rosner, et. al., "Flash Code: Studying Astrophysical Thermonuclear Flashes," Computing in Science and Engineering, 2:33, 2000.
9. J. Cummings and M. Aivazis and R. Samtaney and R. Radovitzky and S. Mauch and D. Meiron, "A Virtual Test Facility for the Simulation of Dynamic Response in Materials," The Journal of Supercomputing 23(1), pp. 39–50, August 2002.
10. C. Hill and C. DeLuca and V. Balaji and M. Suarez and A. da Silva, "The Architecture of the Earth System Modeling Framework," Computing in Science and Engineering, 6(1), Januaray/February 2004.
11. W. Cabot and A. Cook and C. Crabb, "Large-Scale Simulations with Miranda on BlueGene/L," Presentation from BlueGene/L Workshop, Reno, October 2003.
12. A. Wissinsk and R. Hornung and S. Kohn and S. Smith and N. Elliott, "Large Scale Parallel Structured AMR Calculations using the SAMRAI Framework," Proc. SC'2001 Conference, ACM/IEEE, 2001.
13. Y. Zhang and S. Chandra and S. Hariri and M. Parashar, "Autonomic Proactive Runtime Partitioning Strategies for SAMR Applications," Proc. NSF Next Generation Systems Program Workshop, IEEE/ACM 18th IPDPS Conference, April 2004.