

# Performance Analysis, Modeling and Enhancement of Sandia's Integrated TIGER Series (ITS) Coupled Electron/Photon Monte Carlo Transport Code

Mahesh Rajan, Brian Franke, Robert Benner, Ron Kensek and Thomas Laub  
Sandia National Laboratories

**Abstract**— ITS is a powerful and user-friendly software package permitting state-of-the-art Monte Carlo solution of linear time-independent coupled electron/photon radiation transport problems, with or without the presence of macroscopic electric and magnetic fields of arbitrary spatial dependence. As one of a few Sandia applications that are targeted for capability class machines like the ASC Red Storm, we have studied extensively the performance of this application using thousands of processors. We have successfully constructed a performance model and verified the model against measurements on a variety of Sandia compute platforms. Use of tools like VAMPIR and PAPI in performance analysis and modeling is discussed. The original algorithm for computing the statistical quantities after each batch of Monte-Carlo computations has been modified to yield improved parallel scaling. Models of alternate message passing algorithms are investigated and validated against measurements on the Red Storm.

**Index Terms**—Performance modeling, performance analysis, Monte Carlo Radiation Transport

## I. INTRODUCTION

THE INTEGRATED TIGER SERIES (ITS) code is an evolving Monte Carlo radiation transport code that has been used extensively in weapon-effect simulator design and analysis, radiation dosimetry, radiation effect studies and medical physics research. Many individuals from the DOE labs and NIST have been involved over the years in the development and enhancement of ITS [1]. The different features/sections of the code in ITS: TIGER, MITS, CEPXS, XGEN etc., are applied to an analysis under investigation through the selection of appropriate pre-processor directives when the code is built. Physical rigor for the analysis is provided by employing accurate cross sections, sampling distributions, and physical models for describing the production and transport of the electron/photon cascade from 1.0 GeV down to 1.0 keV. The ITS code is capable of analyzing particle transport

through both combinatorial geometry models and CAD models. It also has been significantly enhanced to permit adjoint transport calculations.

For the purposes of this paper we have analyzed the performance using as input, data from a real satellite model. The physical problem solved takes advantage of the MITS mutli-group/continuous energy electron-photon Monte Carlo transport code's capability to address realistic three-dimensional adjoint computations. The adjoint transport method is a powerful technique for simulating applications where the knowledge of the particle flux is only required for a restricted region of the phase space, but where this knowledge is required for source parameters spanning a large region of phase space. The run times for simulations for a complex combinatorial geometry model using conventional, or forward, transport are prohibitive and hence the adjoint calculations used in our satellite model [2].

Our performance analysis of the ITS code was initially spurred by the JASONS and NAS review of the ASC programs to assess mapping of a set of DOE applications to architectures. Another reason for this investigation is because of the large percentage of compute resources ITS code users had consumed in the previous years and anticipated similar usage in the future. Recently we investigated the scaling characteristics of ITS to tens of thousands of processors. Execution time measurements have been obtained on various platforms at SNL; ASCI Red, VPLANT (2.4 GHz Xeon cluster with Myrinet), ICC (institutional cluster: 3 GHz Xeon with Myrinet), CPLANT (Alpha cluster with Myrinet) and more recently on the Red Storm. Our performance model attempts to follow a similar approach to that expounded by Kerbyson, et.al [4] and in fact follows closely the model presented by Mathis, Kerbyson and Hoise [5] in their analysis of the MCNP particle transport code. The model develops an analytical expression for the major portions of the execution time, namely, computation, communication and I/O. At the present time our expression for the compute time is obtained by curve-fitting the plot of the measured execution time vs. the number of histories. For the communication time model we focused on the communications at the end of each batch of computations assigned to the processors. This was accomplished with the VAMPIR tracing tool to obtain the message sizes and messaging patterns and later correlated to the MPI calls in the code. For the compute platform

Manuscript received August 26, 2005. This work was supported in part by the U.S. Department of Energy, SNL CSRF project

Authors are with the Sandia National Laboratories, P.O.Box 5800, Albuquerque, NM, 87185 (Contact phone: 505-284-5063; fax: 505-844-2067; e-mail: mrajan@sandia.gov).

Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States National Nuclear Security Administration and the Department of Energy under contract DE-AC04-94AL85000.

communication characteristics such as bandwidth and latency, a set of simple benchmarks were run. The I/O time typically been a fraction less than 3% and at the present no model has been investigated.

In this paper we present the details of the model and compare measured performance against the model for different architectures. We also present results of tests on the new ASC Red Storm and use the model to predict performance on it to 10,000 processors. The ITS code has been recently enhanced to introduce Fortran 90/95 features and in the process it has also implemented changes in data structures that would improve performance. As will be seen from the scaling plots presented below, ITS can suffer a scaling performance penalty depending on how the history computations are split among the participating processors and the frequency with which the statistical tally of the computations are assembled by the master process. This performance penalty, due to communication cost incurred in the many-to-one communication at the end of each batch of computation, has been remedied by a modified algorithm. We have investigated a simple implementation using MPI collective communication calls and measured its performance on the Red Storm. The MPI collective communication calls (based on MPICH 2.0) are implemented using a binary reduction algorithm and executed using MPI point-to-point operations [6]. We also investigate potential improvements in performance that could result from using Rebenseifner's [7] algorithm that promises better efficiency when message lengths of varying sizes are communicated.

Finally we also present some PAPI hardware performance counter results. We are hoping to use these results in tuning single processor performance and in understanding memory access patterns of the code to evolve a single processor performance model. We believe our analysis and code improvements will enhance the productivity of ITS code users and permit effective usage of new ASC capability class machines like the Red Storm. There is a strong impetus from DOE to push the production use of such systems to utilize upwards of 4000 processors on a regular basis.

## II. ITS CODE FLOW AND SATELLITE COMPUTATIONS

Description of the code and details on using ITS can be found in ITS User Guide [7]. In this paper, we present a broad outline of the computation phases and the parallelization strategy used in ITS. The interaction between particles and the physical geometry under consideration is analyzed by tracking particle trajectories through the geometry. Typically there is a linear relationship between the number of particle trajectories tracked and the execution time once the number of histories exceeds a few thousand. The ability to conduct independent computation of particle histories on a geometry replicated on every processor, combined with the fact that the statistical uncertainty in any Monte Carlo result decreases as the square

root of the number of histories makes it ideally suited for parallel computations. While there exists particle transport codes in which the geometry may also be distributed, like in LLNL's Mercury Code[8], our analysis is simplified because ITS replicates the geometry on all the processors. The basic computation steps are as follows: (1) the master processor reads geometry, problem input and broadcasts to all the worker processes; the number of histories to be computed is divided among the processors and these in turn may be further divided into batches to facilitate break up of computations for convenience in getting intermediate results and keeping the time of computations between batches correlated to the restart dumps, (2) the worker processes perform the Monte Carlo computations, and (3) the master performs Monte Carlo computation and receives and tallies the data after each batch of computations for statistical calculations and outputs the data. Thus based on how the problem is set up we could have many batches of computations assigned to each processor and after completion of each batch there being a communication operation of the tally data from the workers to the master. So although there is no communication at all between the worker processors, there is many-to-one communication which, based on the algorithm used, has significant impact on the parallel scaling characteristics.

As mentioned in the introduction this analysis was undertaken with a satellite combinatorial geometry model. The calculations performed for this work were adjoint point estimation of KERMA (Kinetic Energy Released per unit MA<sub>ss</sub>). This simulation is performed by modeling only photon physics. The adjoint solution allows numerous radiation sources to be assessed for a single detector. In this case, the detector was measuring energy deposition at a point inside of an electronics box. The forward sources assessed were infinite plane-wave sources of photons with a uniform energy distribution between 1 and 50 keV. The average response to these sources was calculated for 472 angular bins of approximately equal solid angle over all  $4\pi$  possible incident directions. Figure 1 illustrates the dosage computations where the pixels are angular bins of the source directions and the levels are dose values at the same point on the object.

The combinatorial-geometry (CG) model of the satellite comprised of a total of more than 600 CG bodies. Portions of the satellite are approximated as reduced-density materials distributed over regions that in reality are a combination of void and intricate geometries, but the model is among the most complicated CG models that we have available. Within the satellite, an electronics box is more accurately modeled with more than 200 CG bodies. For purposes of this investigation we focused on a scaled-speedup analysis (weak-scaling) with 3.2 million histories assigned to each processor.

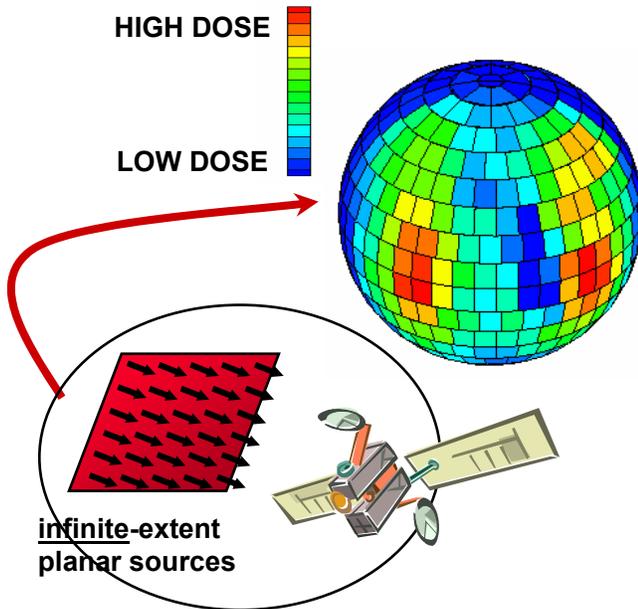


Figure 1. Adjoint calculations for the satellite model assesses the directions of vulnerability

### III. PERFORMANCE MODEL

The performance model for ITS parallel computation consists of two parts: the computation time model and the communication time model. There is a certain amount of time spent in I/O, but for typical calculations running several hours this turns out to be a negligible fraction of the total run time. Following the approach taken in Ref. [4, 5] the computation model is related to an intrinsic parameter of the computation, which here is simply the number of particle histories,  $N_{ph}$ . We measure this quantity using a single processor on each system. The computation model also depends on the geometry and therefore such measurements must be carried out for every geometry. For the satellite geometry the relationship between single processor execution time and number of histories is a linear relationship as shown in Figure 2, for the VPLANT Cluster.

From such measurements the computation time can be calculated as:

$$T_{compute} \equiv N_{ph} * T_{hist} / P$$

where  $T_{hist}$ , time for a single particle history, is obtained from the slope of the line in Figure 2.  $P$  is the number of MPI processes taken to be equal to the number of processors.

The communication time can be modeled as:

$$T_{communication} = T_{setup} + T_{tally}$$

$T_{tally}$  refers to the worker-master communication time after each batch of Monte Carlo computations to gather the statistics for all the histories computed. This gather operation

in ITS was structured as a many-to-one communication resulting in the communication time being directly proportional to the number of processes used in the computation.

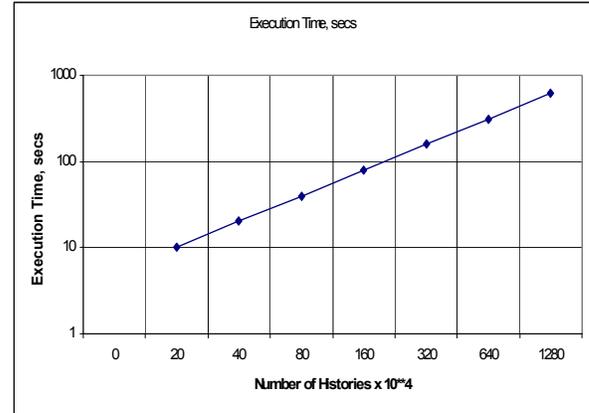


Figure 2. Measurement of single processor execution time on SNL's VPLANT Cluster

ITS has as the option of assigning all the history computations to the worker processes leaving the master to coordinate the computations, this option was not used in our scaling studies. Therefore p-1 processes serially send the tally data to the master. The amount of data that is gathered by the master processor is related to the output of physical quantities requested by the analyst via the input deck and the size of these data elements depend again on the problem under investigation as they are related to the size of the physical arrays. In ITS the size of most of the important physical arrays are controlled through an include file in static fashion. Initially in building the model, we used the VAMPIR performance analysis tool to obtain the details of the communication and later correlated it with MPI calls in the code.

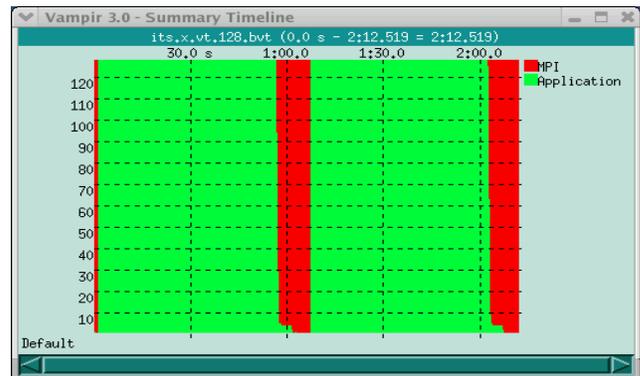


Figure 3. A VAMPIR summary timeline on SNL's ICC cluster

One could zoom into the red region in Figure 3 to gather the communication details. For the satellite model under consideration the details of the messages between a typical

worker and the master is captured in Figure 4.

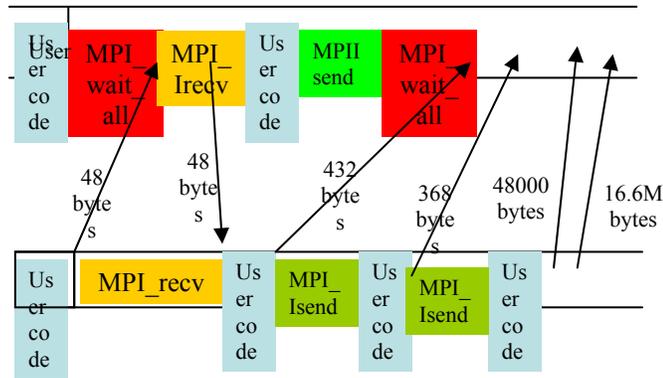


Figure 4. Worker-Master tally data communication calls

As seen from Figure 3, the communication time for setup is a small fraction. Therefore in our model we have chosen to ignore it although one could build a model for it following a similar procedure. So the analytical expression for the communication time is:

$$T_{comm} = \{ 2 * T_{48} + T_{48000} + T_{432} + T_{16M} + T_{368} \} * num\_batches * (p-1)$$

Each processor at the end of a batch of computation sends to the master a data structure coded in an 8 byte integer array called batch\_info(6) that signals the master on completion of its task and some book keeping information such as the execution time for that batch. The master receives the data and using the same data structure instructs the worker process to transmit the tally data for the physical quantities ( such as TAU, ALINE, etc in the code). The 432, 368, 48000, and 16M byte messages corresponding to the various physical quantities are then received by the master from each worker processor one after another in a serial fashion, from the p-1 processes.

To compute the different terms in the communication time expression above, we measure the point-to-point communication time as a function of the message size using simple MPI benchmarks. This measured information, shown in Table 1, on bandwidth and latency are used in the calculation of the communication time as a function of the number of processors and batches. The total execution time is readily tabulated against the number of processors, from which parallel efficiency is calculated. Figure 5 compares the parallel efficiency obtained with our model against performance measurements on ASCI Red (janus), CPLANT, VPLANT and SNL’s institutional cluster – ICC-LIBERTY.

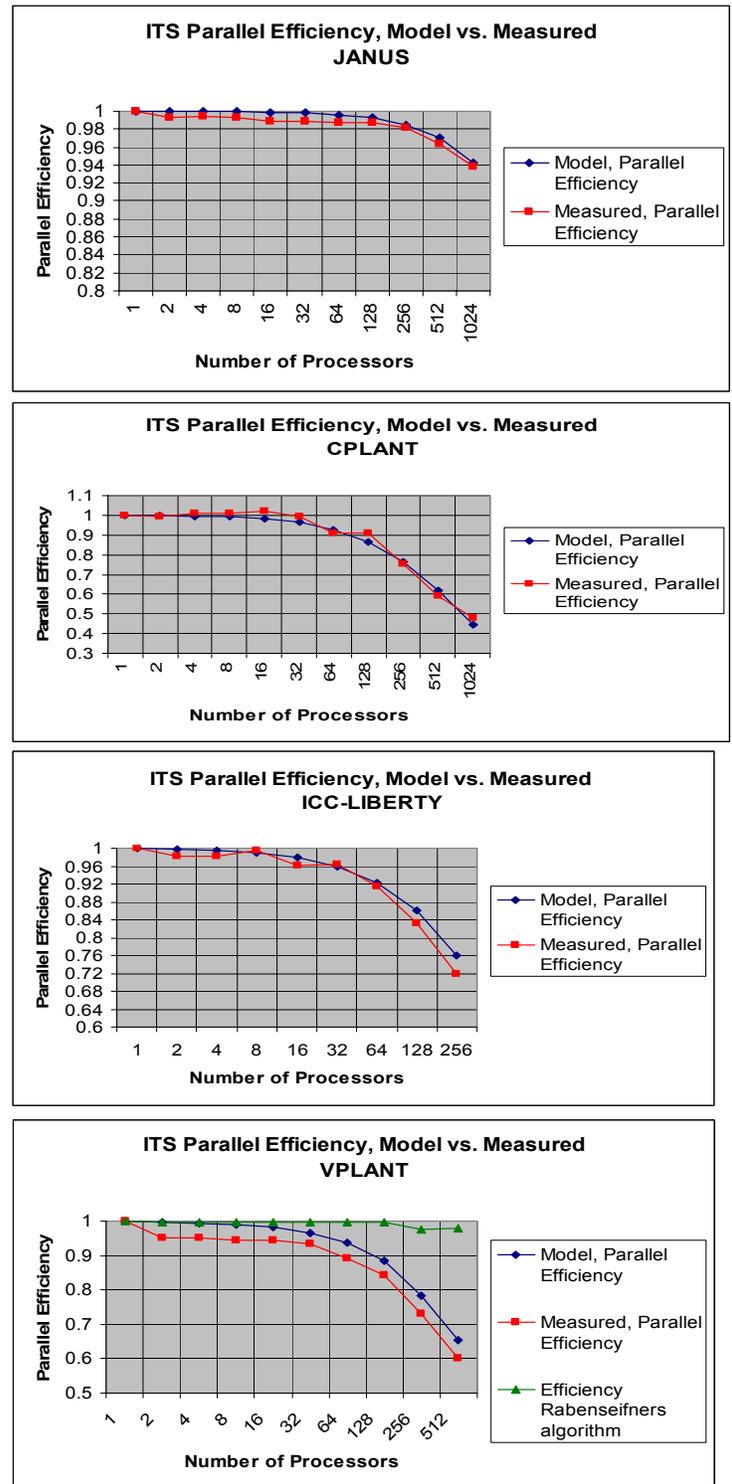


Figure 5. Comparison of parallel efficiency; Model vs. Measurements on ASCI Red, CPLANT, ICC, and VPLANT.

An objective of our performance modeling effort is not only to understand performance characteristics of applications, but

to implement changes in algorithm that would improve performance. Towards this end alternate implementations of the gather phase of the computation using some binary reduction algorithms were investigated. Thakur and Gropp[7] discuss efficiency of various algorithms for MPI collective operations which are now part of the MPICH release (version 1.2.6 and higher). Before we actually implemented a change we analyzed the potential improvement in parallel efficiency using Rebenseifner's [8] algorithm. The expression for the reduction operation for this algorithm is given by:

$$T_{communications} = 2 \ln(p)\alpha + 2(p-1)\beta n + (p-1)n\gamma / p$$

Here  $p$  is the number of processors,  $\alpha$  is the latency,  $\beta$  is the message transfer time per byte and  $\gamma$  is the local reduction time per byte and  $n$  is the message length. The improved scaling that would result from implementing this algorithm is seen in Figure 5 for the data corresponding to the VPLANT cluster.

Table 1. System parameters used in the computation model and comparison of parallel overhead at 512 processors

System	Pt-to-pt BW MB/s	Pt-to-pt Latency, usec	Comp time, secs	Comm time, secs	Overhead, Parallel Efficiency, f & (1/1+f)
Red Storm	857	7	246.92	19.44	0.078, 0.927
Janus	330	18	1673	53.20	0.03, 0.97
ICC	245	6.8	108	69	0.63, 0.61
VPLANT	209	7.9	156	83	0.53, 0.65
CPLANT	76	40	334	237	0.70, 0.58

Recently we have applied the performance model to the new ASC Red Storm[8] and used the model to predict parallel efficiency to 10000 processors. ITS code has also recently undergone major rewrite to take advantage of the new FORTRAN 90/95 features. This rewrite has created new data structures that has facilitated investigation of the alternate message passing scheme to replace the  $O(p)$  dependent communication scheme to  $O(\ln(p))$  communication algorithm. In Red Storm the currently installed MPI software based on MPICH takes advantage of the improved collective communication calls. So the code was modified to simply use MPI collective operations. The results of this improvement can be seen with measured parallel efficiency curve with the 'new code' in Figure 6. If the calls to MPICH collective operations are replaced with Rabenseifner's MPI collective algorithm we should see further improvement in performance as it promises to be more efficient with both long and short messages. Measurements and models with alternate algorithms will be included in the final presentation .

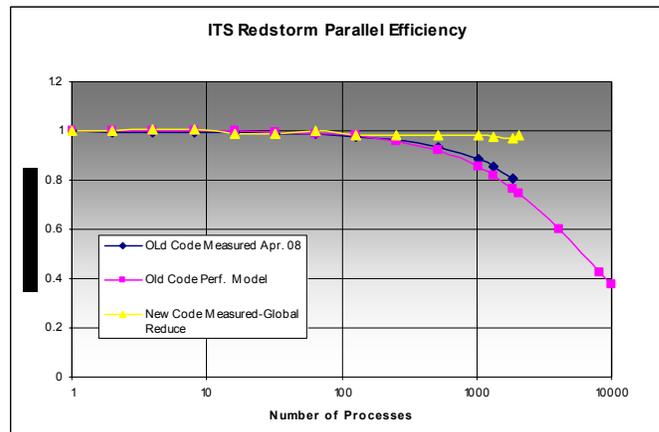


Figure 6. Red Storm; model and measured data showing improved performance with the 'new code'

#### IV. PAPI ANALYSIS AND SINGLE PROCESSOR TUNING

Using the new ITS F90/F95 code we have instrumented the code to get hardware performance counter information using PAPI. PAPI analysis shows that ITS has only a small percentage of the instructions that are floating point operation. Moreover, the nature of the particle trajectory computations has a lot of code branching. This seems to suggest that proper speculative execution with profile based compiler optimization on the Xeon, Em64T and Opteron processors might yield some significant single processor performance improvement. This data is presently being collected and analyzed and the final presentation will include further discussion on this topic.

#### V. CONCLUSION

We have successfully built and used an analytical performance model for ITS to understand its performance on all of Sandia's major computational resources. Our approach follows closely the approach expounded by LANL's PAL team. The performance model and analysis has helped us to identify bottlenecks in communication which were limiting the scalability of this application. Along with other code improvements we have modified the communication algorithm permitting good scaling of this application to  $O(10K)$  processors.

#### REFERENCES

- [1] J.A. Halbleib, R.P. Kensek, T.A. Mehlhorn, G.D. Valdez, S.M. Seltzer, and M.J. Berger, "ITS Version 3.0: The integrated TIGER Series of Coupled Electron/Photon Monte Carlo Transport Codes," Technical Report SAND91-1634, Sandia National Laboratories, 1992.
- [2] R.P. Kensek, et al., "DTRA High Performance Computing for Testable Hardware Initiative Final Report," Sandia National Laboratories, SAND2001-2900 (2001).
- [3] Ang, J.A., *personal communications*, August 2005 .
- [4] Kerbyson, D.J. et al., " Predictive Performance and Scalability of Modeling of a Large-Scale Application", In the Proceeding of the

IEEE/ACM conference on Supercomputing Sc '01, Denver, CO,  
October 2001

- [5] M.M. Mathis, D.J. Kerbyson, A. Hoise, "A performance Model of nondeterministic Particle Transport on Large-Scale Systems" In Future Generation Computer Systems, To Appear 2005, LA-UR 02-7313
- [6] R. Rabenseifner, "Optimization of Collective Reduction Operations", International Conference on Computational Science, June 7-9, Krakow, Poland, LCNS, Springer-Verlag, 2004.
- [7] M. O'Brien, J. Taylor, R. Procassini, "Dynamic Load Balancing of Parallel Monte Carlo Transport Calculations," The Monte Carlo Method: Versatility Unbounded In A Dynamic Computing World, Chattanooga, Tennessee, April 17-21, 2005, on CD-ROM, American Nuclear Society, LaGrange Park, IL (2005)
- [8] Hoise, et.al., "An initial Performance Analysis of the Red Storm Architecture" LANL PAL on-line distribution March 14, 2005