# Architectural Requirements for the Red Storm Computing System

James L. Tomkins and William J. Camp

Sandia National Laboratories

# Architectural Requirements for the Red Storm Computing System

James L. Tomkins
Computer & Software Systems

William J. Camp
Computation, Computers, Information & Mathematics

Sandia National Laboratories
P.O. Box 5800, MS-1109
Albuquerque, NM 87185-1109

## Abstract

This report is based on the Statement of Work (SOW) describing the various requirements for delivering a new supercomputer system to Sandia National Laboratories (Sandia) as part of the Department of Energy's (DOE) Accelerated Strategic Computing Initiative (ASCI) program. This system is named Red Storm and will be a distributed memory, massively parallel processor (MPP) machine built primarily out of commodity parts. The requirements presented here distill extensive architectural and design experience accumulated over a decade and a half of research, development and production operation of similar machines at Sandia. Red Storm will have an unusually high bandwidth, low latency interconnect, specially designed hardware and software reliability features, a light weight kernel compute node operating system and the ability to rapidly switch major sections of the machine between classified and unclassified computing environments. Particular attention has been paid to architectural balance in the design of Red Storm, and it is therefore expected to achieve an atypically high fraction of it's peak speed of 41 TeraOPS on real scientific computing applications. In addition, Red Storm is designed to be upgradeable to many times this initial peak capability while still retaining appropriate balance in key design dimensions. Installation of the Red Storm computer system at Sandia's New Mexico site is planned for 2004, and it is expected that the system will be operated for a minimum of five years following installation.

# Acknowledgements

4

**(This plane left intentionally blank)**

# Architectural Requirements for the Red Storm Computing System

James L. Tomkins
Computer & Software Systems

William J. Camp
Computation, Computers, Information & Mathematics

Sandia National Laboratories
P.O. Box 5800, MS-1109
Albuquerque, NM 87185-1109

## Introduction

This report is based on the Statement of Work (SOW) describing the various requirements for delivering a new supercomputer system to Sandia National Laboratories (Sandia) as part of the Department of Energy's (DOE) Accelerated Strategic Computing Initiative (ASCI) program (which includes Sandia, Los Alamos, and Lawrence Livermore National Laboratories). This system is to be named Red Storm and will be referred to throughout this SOW as either the Red Storm or the 40 TeraOPS computer system. Installation of the Red Storm computer system at Sandia's New Mexico site is planned for 2004, and it is expected that the system will be operated for a minimum of five years following installation.

The supercomputer will be built from components that, in the 2003 - 2004 time frame, will be Commercial Off The Shelf (COTS) as much as reasonably possible while still meeting the requirements of this SOW. Red Storm will represent a major increase in computing capability for both Sandia and the ASCI program.

The statement of work contains two parts corresponding to the two major tasks required of prospective vendors. Task I is for the development, design, manufacture, test, assemblage, installation and maintenance of the Red Storm system and Task II is for the program management and support of the Red Storm system. An optional upgrade to the Red Storm system to 120 TeraOPS or greater shall be planned for as part of Task I.

## System Performance Background

Historically the performance of supercomputers has been measured in a number of ways including by peak Floating-point Operations Per Second (OPS), by simple benchmarks such as

MPLINPACK, and by complex physical simulations. The best current supercomputers have achieved 70-75% of peak performance on the MPLINPACK benchmark. However, for many complex simulation codes the performance is only 10-20% of peak for a single processor and can be as low as a one or two percent when parallel efficiency is considered. The performance, as measured against peak, for complex simulation codes has been declining in recent supercomputer generations. Regrettably, this trend seems to be continuing in the newest supercomputers.

The two most significant areas of computer hardware design that have contributed to this trend of reduced performance relative to peak are in the machine interconnect and node memory system. Interconnect hardware development has severely lagged the pace of increasing processor performance. The shift from tightly coupled Massively Parallel Processor (MPP) designs such as the Intel ASCI Red and Cray T3E designs to clusters that use I/O buses for interconnect connections has resulted in not only a relative reduction in interconnect performance but an absolute reduction. At the same time, processor performance has been increasing rapidly. This combination has resulted in growing performance imbalance in large parallel computer systems. Also, the size of machines in terms of the number of processors has been increasing, putting even more stress on interconnect performance. The result has been poor application scalability compared to that achieved on earlier generations of tightly coupled MPPs and poor overall efficiency of computer systems.

Node memory system performance for micro-processor based computers has been increasing but at a much slower rate than peak CPU performance. This is true both in terms of bandwidth and latency. Recent developments such as dual data rate SDRAM and RAMBUS memory may have slowed this trend and, for some node designs, may have even reversed the trend temporarily. Larger and faster L2 caches have helped to mitigate the growing imbalance in memory system performance, however, not all application codes can make good use of cache. Another factor that influences node memory system performance is the size of the cache line. For applications that access memory in a non-uniform pattern much of the memory bandwidth can go to waste. Node memory system performance is very important to overall performance for scientific applications and the only real solution is to develop memory systems that can keep up with the CPUs.

There are several areas in which system software development or the lack thereof is having a negative impact on the performance of complex simulation codes. The major issue is the scalability of the operating system and operating system services such as job loading, internal communication, network communication, file management, and file I/O.

A second major area in which software development is not keeping up with processor performance increases is in compiler technology. While the peak performance potential of processors has and continues to grow rapidly with each new generation of processors, compiler developers are struggling to take advantage of the hardware performance increases. The problem of translating user code into efficient machine code is becoming more and more difficult. In addition, the trend in language development has recently been toward language constructs which encourage inefficiencies by stressing memory systems through indirect addressing associated with complex data structures, large numbers of compiler generated temporary variables, compiler generated memory coping, and by encouraging the use of array operations which result in little or no cache reuse.

The performance goal for the Red Storm system is a minimum of at least a factor of 7 improvement in the average performance for a suite of ASCI simulation codes as compared to the current ASCI Red system. These codes will encompass a wide variety of applications including shock physics, radiation transport, materials aging and design, computational fluid dynamics, structural dynamics, and others. They will also include structured and unstructured grids, explicit solvers, and sparse matrix solvers, and they must scale to the full size of the system.

## Reliability, Availability, and Serviceability (RAS) Considerations

ASCI scientific and engineering calculations need to run for a 100 or more hours on a large fraction of or even in some cases the whole machine. When a failure occurs that causes a calculation to be interrupted all of the computing since the last completed restart dump is lost. For example, a problem running on 10,000 processors that has a two bit error in the memory associated with a single processor will fail on that processor with an unrecoverable error. The other 9,999 processors will be halted and the job killed. In general it is not possible to recover the calculation by only redoing the part of the calculation that failed because to do so would require recreating all of the communication, from the time of the last restart dump, that was sent to the processor where the failure occurred. Keeping all of the communication traffic for a problem running on 10,000 processors would rapidly exceed any reasonably possible amount of disk storage and would swamp the machines I/O system because it would have to be kept for all 10,000 processors. Currently the only choice is to go back to the previous restart dump and restart the full calculation, losing all of the work since the last restart dump. Frequently performing restart dumps would decrease the amount of lost work when an interrupt occurs, but would also increase the overall run time for the calculation by adding overhead for writing the restart dumps. It also stresses the computer systems disk I/O capabilities as it increases the total amount of I/O needed (more restart dumps) to complete the calculation. Because of the types of applications, scientific and engineering, that are being run on the ASCI machines, the key criteria for measuring reliability is Mean Time Between Interrupts (MTBI). System availability is of secondary importance. In fact, it is possible to have a machine with high availability that is not useful for ASCI problems because its MTBI is too short.

## Classified and Unclassified Computing Considerations

Sandia and the ASCI program need to be able to do both large scale classified and unclassified computing. At Sandia, almost all of the application code development is being done in the unclassified partition while much of the calculation work load is classified and must be done in the classified partition. In addition, ASCI has a requirement to provide large scale, unclassified computing access to its Alliance Program partners. The only practical way to meet these requirements is through switching a large fraction of the computing capability between the classified and unclassified computing partitions.

# Task I: System Requirements - Develop, Design, Manufacture, Test, Assemble, Install, and Maintain the Red Storm MPP Supercomputer System

Sandia has a long history in MPP computing starting with the first 1024 processor nCUBE 10 computer system that was installed at Sandia in 1987. Since then Sandia has had two 1024 processor second generation nCUBE 2 MPPs, a 16K processor Thinking machines CM2, a 3600+ processor Intel Paragon and currently the Intel ASCI Red machine with over 9500 processors and over 2500 processors in Cplant clusters. Over this period computing technology has changed dramatically, however, the basic characteristics that made the first nCUBE 10 a highly scalable machine and ASCI Red highly scalable are the same. These are a communication network providing a high computation to communication ratio, a highly scalable operating system and system software, an integrated system design, and a level of reliability that allows for meaningful work to be accomplished between interrupts. The requirements that follow were developed based on our extensive experience with MPPs and are designed to produce a computer system that will meet our needs.

In developing this set of requirements Sandia desires to leave, as much as possible, the technical methodology of achieving our computational performance goals to the contractor. However, our experience with developmental computer systems has led us to believe that it is absolutely necessary for us to set minimum requirements for performance and reliability if we are to hope to achieve our goals. Also, because we have some special needs, for example to be able to switch a major portion of the machine between classified and unclassified computing, it is necessary for us to have more specific architectural design requirements than otherwise would be required.

The following sections describe the minimum system requirements for the Red Storm computer system.

## Red Storm System Requirements

These system requirements were developed based on Sandia's long and successful experience at applying MPP computers to real application problems in high performance technical computing. These requirements reflect our goal of achieving at least a factor of 7 improvement in the average performance on a suite of ASCI codes as compared to the current ASCI Red system. They also reflect Sandia's and the ASCI programs special needs for security and flexibility in using the system in both the unclassified and classified environments.

Red Storm system requirements are presented in Tables 1 - 7. They include hardware architecture and performance requirements: system software functionality and performance requirements; system reliability, availability, and serviceability requirements; and system security requirements. Maintenance requirements for the Red Storm system are also included under Task I after Table 7.

The contractor will develop, design, manufacture, test, assemble, install and maintain a 40 TeraOPS1 MPP supercomputer system (Red Storm system) that satisfies all of the requirements

in Tables 1 - 7 below. The Red Storm system shall be able to achieve a sustained 28 TeraOPS on the MP LINPACK benchmark, shall have at least 10 TBytes[2] of accessible RAM memory on the compute nodes, and shall have at least 240 TBytes[2] of formatted user disk storage.

## System Architecture Requirements

Sandia's long history and experience in large scale MPP computing together with the special needs for both classified and unclassified computing drive the system architecture requirements given in Table 1. These requirements present a computer architecture that is similar to the ASCI Red architecture that has proven to be highly successful for ASCI scale computing. They also address the special security needs of Sandia and the ASCI program. The system layout is given in Figure 1. The parallel disk storage is not shown but it is connected to the service and I/O nodes which are shown. [In the figure, two of the three hardware partitions are shown.] There are four rows of cabinets with a total of 35 cabinets in each row. The classified service and I/O partition is shown on the left side and the unclassified service and I/O partition is shown on the right side. These partitions are two cabinets wide and across all four rows. Separating the compute partition from the service and I/O partitions is a disconnect cabinet in each row at each end of the machine. Within the compute partition there are 27 CPU cabinets and 2 disconnect cabinets for each of the four rows. For each row, 7cabinets are shown as classified, 13 cabinets as switchable, and 7 cabinets as unclassified. These cabinets are all assumed to be two feet by four feet and an isle width of four feet is assumed. All I/O between the compute partition and disks and external networks goes through the service and IO nodes at each end of the machine. The RAS and system management partition, except for the system management workstations which are not shown in the figure, are contained in the CPU cabinets.

---

1. For the purposes of this SOW performance requirements are given in decimal numbers. 40 TeraOPS is $4.0 \times 10^{13}$ floating-point operations per second.
2. For the purposes of this SOW all storage and memory sizes are based on powers of two. A KByte is 1,024 Bytes (210), a MByte is 1,048,576 Bytes (220), a GByte is 1,073,741,824 Bytes (230), and a TByte is 1,099,511,627,776 Bytes (240).

Figure 1: Red Storm System Layout

# Table 1:
## System Architecture Requirements for the Red Storm Computer System

| Description |
| --- |
| 1.1 The system shall be a tightly coupled MPP designed to be a single system. |
| 1.2 The system shall have a distributed memory Multiple Instruction Multiple Data (MIMD) architecture. |
| 1.3 The Red Storm primary internal communication network shall be a fully connected, true 3-Dimensional mesh topology. The compute node mesh topology shall be 27 X 16 X 24 (x, y, z) based on the minimum number of processors in requirement 1.4 below. If the contractor chooses to use a torus for any direction (x, y, or z) of the Red Storm compute node mesh the full torus connectivity shall be provided for all Red/Black switch configurations listed in 1.9 below. |
| 1.4 There shall be at least 10,368 AMD Opteron (Sledgehammer) compute node processors in the system. The processors will be configured in 108 compute node cabinets which are arranged in 4 rows of 27 cabinets. Each compute node cabinet will contain 96 processors. |
| 1.5 Aggregate system RAM memory for the compute nodes shall be at least 10.0 TBytes. The compute node system RAM memory shall be divided among the compute node processors to provide each compute node processor with an equal amount of local RAM memory. |
| 1.6 There shall be a minimum of 240 TBytes of user formatted disk space that is split into two separate disk systems, one classified and one unclassified, each of 120 TBytes of formatted user disk system storage. |
| 1.7 The Red Storm computer system shall have three functional hardware partitions; 1) Service and I/O, 2) compute, and 3) RAS and system management. |
| 1.8 The Red Storm computer system shall be air cooled. The system cabinets for compute nodes and service and I/O nodes shall be spaced 24 inches on center and be 24 inches wide. This requirement is derived from the need to match the computer room raised floor layout which has support on 24 inch centers. |

# Table 1:
## System Architecture Requirements for the Red Storm Computer System – Continued

| Description |
| --- |
| 1.9 The system shall support Red/Black switching within the compute partition and between the compute partition and service and I/O partitions. The compute partition shall have three sections; 1) unclassified (~1/4), 2) switchable (~1/2), and 3) classified (~1/4). Normal configurations of these three sections shall be ~1/4 unclassified and ~3/4 classified and ~3/4 unclassified and ~1/4 classified. The compute partition (all three sections) shall also be able to be configured as a single unclassified or classified partition. For all of the above possible system configurations, the full connectivity of the 3-D mesh shall be maintained. User access to the classified and unclassified service and I/O nodes and user disk space shall be maintained for all possible system configurations. Connectors used to support the Red/Black functionality of the system shall be designed to provide a minimum of 5,000 connect/disconnect cycles without failure. Reconfiguration of the system from any of the above configurations to any other of the above configurations shall be completed in less than 1 hour (see requirement 5.3 in Table 5). |
| 1.10 There shall be two service and I/O node partitions, one for classified computing and one for unclassified computing. Each service and/or I/O node shall be connected to the primary internal communication network with at least one bi-directional link operating at the full internal communication network link bandwidth. These links shall be designed to provide a minimum of 5,000 connect/disconnect cycles without failure. Each service and I/O node partition shall have 256 AMD Opteron (Sledgehammer) processors configured in 8 cabinets (16 cabinets total) in a processor topology of 2 X 8 X 16 (x, y, z). |
| 1.11 The operating system and system software functionality shall be partitioned to match the hardware partitioning. Service and I/O nodes shall have a full UNIX or UNIX-derivative operating system. Compute nodes shall have a light weight kernel operating system. The RAS and system management system shall have a real-time or real-time like operating system and for the system management workstations a UNIX or UNIX-derivative operating system. |
| 1.12 Source based routing is the preferred method of routing for the primary internal communication network. A system based on router tables may be used provided it is scalable to the maximum size that an upgraded Red Storm system might grow to. **Connection based protocols shall not be used.** The amount of system RAM memory used to provide buffers for message passing shall be independent of the number of compute nodes in the system or of the number of compute nodes assigned to a particular user job. |
| 1.13 The system shall have two disk storage systems, one for classified computing and one for unclassified computing. All system disk storage shall be connected to service and I/O nodes only. Compute nodes shall have no local disk or user writable non-volatile memory. |

# Table 1:
## System Architecture Requirements for the Red Storm Computer System – Continued

| Description |
| --- |
| 1.14 The full Red Storm system shall be less than 8 feet in height and require less than 6,000 square feet of floor space (gross). The machine including service and I/O nodes, compute nodes, system management workstations, the disk storage systems in 1.6 above, and all space needed for access and maintenance of the machine are to be included in the 6,000 square feet limit. The Red Storm physical layout shall fit within the space constraints of room 230 of Sandia's computer annex in Building 880 at Sandia National Laboratories in Albuquerque, New Mexico.<br><br>1.15 The total power required to run the Red Storm system shall be less than 3.5 MW. The cooling required shall be less than 3.5 MW.<br><br>1.16 The system shall be designed to facilitate an upgrade from 40 TeraOPS peak to a minimum of 120 TeraOPS peak through a processor upgrade and additional cabinets. The interconnect shall remain and the total floor space requirement for the upgraded system shall be less than 10,000 square feet. The total power and cooling required for the upgraded system shall be less than 7 MW. An unpriced option for an upgrade to 120 TeraOPS shall be part of the contract for the Red Storm system. |

A high performance 3-D mesh topology is specified as an architectural requirement because it is a good match to our applications, it is highly scalable and it makes Red/Black switching feasible. A torus has not been specified because it requires longer cables, because it will complicate Red/Black switching if it is used in the direction in which Red/Black switching takes place, and because of the need for all I/O to go through service and I/O nodes at each end of the machine. However, the contractor may use a torus provided that it can meet the other requirements. For example, a 3-D mesh interconnect with one or two of the three dimensions connected in a torus could be built without affecting Red/Black switching.

Functional hardware partitioning of the machine is required to meet the Red/Black switching requirements. It also makes it possible to have denser CPU cabinets in the compute partition than would be possible with general purpose server components. This reduces the overall floor space needed for the system. It also reduces the number of cables and the length of cables needed to make the 3-D mesh because more of the mesh connections can be put on the back-plane and confined to individual cabinets.

Cabinet dimensions for the system are very important to the installation, cable routing, and cooling of the system. Standard raised computer room floors have 24 inch on center support. It is important that the weight of the cabinets be over the support. This allows the center sections of the floor tiles to be perforated for cooling and to have holes cut in them for interconnect cables and other cabling that must go under the floor. Because of strength requirements for the tiles it is

not possible to cut away any significant area of a tile on its edge.

Partitioning of the operating system is driven by the need for overall system scalability and by different functionality requirements for different parts of the system. The Red Storm machine will have a little over 10,000 processors that must work together on a single parallel application code. To achieve the high level of scalability required, the compute node operating system must be driven by the application code rather than managed through a set of processes and demons. Operating systems that have time-sharing, demand paging, sockets, graphical interfaces, and similar high level functionality have demonstrated significant scaling problems. As a result a Light Weight Kernel (LWK) operating system is required for the compute nodes. However, for the user interface, I/O to disks, and external I/O to and from the machine a much higher level of functionality is required and that is why the service and I/O node operating system is specified to be a full UNIX or UNIXderivative operating system.

## Aggregate System Performance Requirements

The aggregate system performance requirements for the Red Storm system are presented in Table 2. The minimum acceptable peak performance of the compute partition is 40 trillion (4.0 X $10^{13}$) 64 bit floating-point operations per second. The performance of service and I/O node processors, processors used to provide the RAS system, processors intended primarily for interconnect functionality, and processors whose primary intended functionality is for non floating-point computation cannot be included in meeting this requirement.

## Table 2:
## Aggregate System Performance Requirements for the Red Storm Computer System

| Description |
| --- |
| 2.1 The system shall have a minimum peak floating-point performance of 40 TeraOPS. Only compute node processors shall be included in the determination of the peak performance for the machine. Service and I/O nodes and RAS system processors are not to be included in determining the 40 TeraOPS minimum. |
| 2.2 The sustained performance on the MP-LINPACK benchmark shall be at least 28 TeraOPS using 64-bit IEEE floating-point arithmetic. This level of performance shall be achieved using the compute partition processors only. The MP-LINPACK benchmark is as defined by Jack Dongara, where only floating-point operations which are actually performed can be counted in the operation count. However, no more than $2(N3)/3$ operations, where N is the dimension of the matrix, may be counted. |

16

# Table 2:
## Aggregate System Performance Requirements for the Red Storm Computer System – Continued

| Description |
| --- |
| 2.3 The peak aggregate memory bandwidth from local node main memory to processors in the compute partition shall be at least 55 TBytes per second. For nodes that have processors that share a memory system, calculation of this memory bandwidth shall be based on dividing the memory system bandwidth by the number of processors that share that bandwidth before multiplying by the number of compute processors in the system. Only compute partition nodes shall be included in this calculation.<br><br>2.4 The maximum MPI latency (hardware and software) for sending a zero byte length message from any user process to any other user process within the compute partition of the Red Storm machine shall be less than 5.0 micro-seconds. The maximum MPI latency for sending a zero byte length message from a user process to a user process on a neighboring node within the 3-D mesh shall be less than 2.0 micro-seconds. These measurements shall be made without having the receiving processor in a polling mode and the latency times must be calculated as ping-pong times divided by 2.<br><br>2.5 The measured minimum bi-directional, bi-section bandwidth shall be 1.5 TBytes/s. Message or packet header bytes and bits required for signaling or error correction shall not be counted in meeting this requirement. The maximum message size for messages used to meet this requirement shall be 1.0 MBytes or less. The minimum bi-section bandwidth measurement shall be determined by the sustained bandwidth through the plane of the 3-D mesh interconnect that has the least bandwidth. Bandwidth provided by nodes that are part of the I/O and service partition shall not be included in this measurement.<br><br>2.6 The minimum sustained bi-directional bandwidth from a node to the Red Storm internal communication network shall be 1.5 Bytes per peak floating-point operation per second of the node. The minimum bandwidth of any single connection from a compute node to the internal communication network shall be 1.5 Bytes per peak floating-point operation per second of a single processor on a compute node. The sustained bandwidth used to meet these requirements shall not include bandwidth needed by the system for address traffic, cache coherency, error correction, or other signaling overhead. |

# Table 2:
## Aggregate System Performance Requirements for the Red Storm Computer System – Continued

| Description |
| --- |

2.7 Each individual link in Red Storm's 3-D mesh primary internal communication network shall have a minimum bi-directional bandwidth of 1.8 Bytes per peak floating-point operation per second of a single compute node processor. Bandwidths used to meet these requirements shall not include bandwidth needed by the system for address traffic, cache coherency, error correction, or other signaling overhead. In addition, each link in the Red Storm 3-D mesh primary communications network shall have at least as much bandwidth as the sum of the bandwidths for all the links from one or more compute nodes that are connected to the interconnect at that switch or router. This means that if one and only one compute node is connected to the 3-D mesh at each router by one link then the six (6) links connecting that router to the 3-D mesh shall each have at least as much bandwidth as the link connecting the compute node to the router. Similarly, if one or more compute nodes are connected to a router by more than one link then each of the six (6) links of the 3-D mesh shall have at least as much bandwidth as the sum of the bandwidths of all of the links connecting the one or more nodes. In other words, each link of the 3-D mesh shall be able to at a minimum handle all of the bandwidth that can be added to and/or removed from the network at each point where compute nodes are connected to the network.

2.8 The system shall be able to simultaneously sustain a minimum of 50 GBytes/s of parallel I/O bandwidth to each parallel disk storage system, classified and unclassified. This bandwidth requirement applies to both writing and reading. To meet this requirement the performance shall be measured from an ASCI application code running on at least 25% of the compute node processors and writing or reading a total of at most 50 TBytes of data. The aggregate reads and writes must exceed the aggregate disk system cache by at least a factor of 10 in size.

2.9 The system shall be able to simultaneously sustain a minimum of 25 GBytes/s of external network bandwidth from both the classified and unclassified service and I/O partitions. This bandwidth shall be sustained using the TCP/IP protocol. The minimum sustainable bandwidth (TCP/IP protocol) of any individual external network interface used to meet the 25 GB/s requirement shall be at least 500 MBytes/s.

2.10 Each service and I/O partition, classified and unclassified, shall be able to support the simultaneous use of the system by a minimum of 30 users running in an interactive mode in addition to 30 batch jobs executing.

2.11 Each service and I/O partition, classified and unclassified, shall have a minimum sustained aggregate bi-directional bandwidth between itself and the compute partition of at least 200 GBytes/s through links to the Red Storm primary internal communication network.

## Table 2:
## Aggregate System Performance Requirements for the Red Storm Computer System - Continued

| Description |
| --- |
| 2.12 The system shall achieve a minimum factor of 7 improvement in average performance measured by execution time on a suite of ASCI application codes over the performance of this same suite of codes on the current ASCI Red System. The suite of codes used for this test will be determined prior to delivery of the Red Storm system to Sandia. A minimum of five different parallel ASCI application codes will be used for this test. The codes shall be run as is except for recompiling.<br><br>2.13 The aggregate sustained bandwidth for all links that make up the compute node 3-D mesh interconnect shall be at least 100 TBytes/s. Only the sustained bandwidth of links that make up the compute node 3-D mesh may be counted in determining this bandwidth. Links that connect compute nodes to the compute node 3-D mesh, links from the compute node 3-D mesh to the I/O and service nodes, and links that connect the I/O and service nodes to each other are not to be included in determining this aggregate bandwidth. The sustained bandwidth shall not include bandwidth needed by the system for address traffic, cache coherency, error correction, or other signaling overhead. |

The requirement for sustained performance on the MP-LINPACK benchmark of 28 trillion ($2.8 \times 10^{13}$) floating-point operations per second is independent of the minimum peak requirement of 40 TeraOPS. In agreeing to meet this requirement the contractor is agreeing to provide additional hardware beyond that necessary to meet the 40 TeraOPS peak performance requirement, should that prove necessary. MP-LINPACK is being used as the primary benchmark for the system because it is well understood, it is a test of the full system, and for a machine of this size it will also be a test of system integration and reliability.

The performance requirement for aggregate local memory bandwidth is driven by the need to improve single node performance relative to peak for real ASCI application codes. It is also driven by scalability considerations since the memory system must also support message-passing traffic. This requirement is independent of the number of processors and of the 40 TeraOPS requirement. The requirements for a high bandwidth, low latency internal communication network are driven by the need for applications to achieve high levels of parallel efficiency while running on the full machine. Poor interconnect performance (hardware and system software) is the biggest reason for poor application scalability and poor performance of large parallel computing systems. An application code that achieves 10% parallel efficiency on 1,000 processors can be sped up by a factor of 5 if its parallel efficiency can be increased to 50% while improving the single node performance of the same application by a factor of ten will only improve the overall performance on 1,000 processors by 9%. In addition, improving the parallel efficiency allows the problem to scale over more processors which further decreases the run time and makes it possible to run larger, more complex problems.

The external network bandwidth requirements are driven by the need to move data on and off of

the machine to and from archival storage, to receive or transmit to other ASCI sites, and for viewing. The calculations that will be performed on this machine will generate many TBytes of output data for a single run. It is very important that the users be able to have access to this data through high performance connections to the machine.

The last requirement in this section addresses the need for real performance from the system. The Red Storm system delivered by the contractor shall provide at least seven (7) times the performance of the ASCI Red system on real applications.

## Compute Node Architecture and Performance Requirements

Requirements for the compute node architecture and performance are presented in Table 3.

### Table 3:
### Compute Node and Back-plane Architecture and Performance Requirements

| Description |
| --- |
| 3.1 Each node in the compute partition shall have at least one and a maximum of 16 compute processors. If there is more than one processor per compute node the processors shall have shared, cache-coherent memory. The compute node processor in Red Storm shall be an AMD Opteron (Sledgehammer) running at 2.0 GHz or greater. All compute nodes shall have an equal number of compute processors. |
| 3.2 Each compute processor shall have a floating-point unit that provides full 64-bit (or greater) IEEE floating-point arithmetic in hardware. Each compute processor shall have an integer processing unit that provides full 64-bit integer arithmetic in hardware. Address registers shall be a minimum of 64-bit length. For the purposes of this SOW, arithmetic is defined as addition, multiplication, and division. |
| 3.3 Each compute processor shall provide user accessible performance counters for the following operations: <br>     1. floating-point adds and multiplies, <br>     2. integer adds and multiplies, <br>     3. load and store operations from cache and from main memory, <br>     4. cache misses. <br> This capability may be provided through hardware or efficient performance monitoring software. |
| 3.4 Each processor shall have a minimum of 1.0 GBytes of uniform access main RAM memory. All RAM memory shall have ECC error correction that is able to detect at least 1 and 2 bit errors and correct single bit errors. For nodes that have more than one compute processor there shall be at least 0.75 GBytes of uniform access main RAM memory per compute processor. |

# Table 3:
## Compute Node and Back-plane Architecture and Performance Requirements - Continued

| Description |
| --- |

3.5 The maximum average latency for a processor to randomly access (fetch or store) 64 bits to or from main RAM memory shall be no more than 325 nano-seconds (325.0e-9s).

3.6 Each node shall have at least one bi-directional connection (Network Interface) to the internal system communication network and not more connections than the number of processors in the node.

3.7 There shall be a minimum sustained memory bandwidth of 2 Bytes per peak floating point operation per second for each processor on each node. The Streams Triad Benchmark (John D. McCalpin) running simultaneously on all processors on a node shall be used to make this measurement.

3.8 Maximum MPI latency for sending a zero byte length message between adjacent nodes shall be less than 2.0 micro-seconds. Maximum MPI latency for sending a zero byte length message between any two nodes in the compute partition of the machine shall be less than 5.0 micro-seconds. These latency measurements shall be made for messages between user processes and without having the receiving processor in a polling mode. The latency requirements are for ping-pong time divided by 2. The measurements are to be made from an application code running in an environment in which other codes are also running on the machine.

3.9 Each interface to the 3-D mesh primary communications network shall be able to sustain a minimum of 1.5 Bytes of bi-directional bandwidth per peak floating-point operation per second of a compute node. This bandwidth requirement shall be met without including any bandwidth required by the system for address traffic, cache coherency, error correction, or other signaling overhead.

3.10 Each switch in the 3-D mesh primary communications network shall have a minimum of 7 bi-directional ports (6 ports for the 3-D mesh and 1 port for the node connections) that have a minimum sustained bandwidth of 1.8 Bytes per peak floating-point operation per second of a compute node processor. This bandwidth requirement shall be met without including any bandwidth required by the system for address traffic, cache coherency, error correction, or other signaling overhead.

3.11 Switch chips in the 3-D mesh primary communications network shall include counters for measuring message traffic performance and for counting recoverable and non-recoverable errors.

For the purposes of this SOW a compute node is defined as having one or more compute processors that share cache coherent memory. Compute nodes shall have at least one, and not more than the number of compute processors on the compute node, bi-directional links to the

machine's primary internal communication network.

The total amount of memory specified in Table 1 and the minimum amount of memory per processor specified here in Table 3 must both be met for the requirements to be met. This is true even if the total amount of memory that would be needed to meet requirement 3.4 exceeds the 10.0 TBytes specified in Table 1.

As processor clock speeds have been increasing, the latency of access to memory in terms of CPU clocks has been increasing. This trend has resulted in reduced real performance for memory intensive applications, especially for those applications involving non-stride-one accesses. This trend has been somewhat offset by improvements in size and speed of cache memory. However, the move to unstructured grids and the use of complex data structures in the application codes has added stress to memory system performance due to increased occurrence of non-stride-one memory accesses. Sandia believes that memory latency is important to overall compute node performance. Therefore, as-low-as-possible memory latency is desired.

Overall system performance depends mostly on the parallel efficiency achievable in the system. High parallel efficiency for a broad spectrum of scientific and engineering codes can only be achieved with a very high performance communications network. The 3-D mesh primary communications network (network interface and switch) specifications have been designed to provide a very high performance communication network.

# Service and I/O Node and Disk System Architecture and Performance Requirements

Service and I/O node and disk system architecture and performance requirements are given in Table 4. These requirements are in addition to the overall system requirements given in Tables 1 and 2. The service and I/O nodes will provide the user interface to the compute partition and as such they need to run a full UNIX or UNIX-like operating system with time-sharing, demand paging, sockets, graphical interfaces, compilers, debuggers, accounting tools, performance monitors, etc.

## Table 4:
## Service and I/O Node and Disk System Architecture and Performance Requirements

| Description |
| --- |
| 4.1 Each service and I/O node shall have a connection to the primary communication network of the compute partition. Each of these connections shall be at the full bi-directional link bandwidth of the primary communication network. |
| 4.2 Disk storage shall provide data integrity protection such as that provided by RAID 3 and RAID 5. Each RAID shall have parity protection in addition to its own hot spare disk. (The preferred configuration is 8 + 1 + 1). |

| Description |
| --- |
| 4.3 The minimum sustained transfer rate for reading or writing of a single UNIX File System (UFS) file to any single logical disk subsystem anywhere in the configured system shall be at least 500 MBytes/s. The file size for this read or write must be at least 10 times as large as the disk system cache on a single disk controller. |
| 4.4 The system shall be able to simultaneously sustain a minimum of 50 GBytes/s of parallel I/O bandwidth to each parallel disk storage system, classified and unclassified. This bandwidth requirement applies to both writing and reading. To meet this requirement the performance shall be measured from an ASCI application code running on at least 25% of the compute node processors for both writing and reading a total of at most 50 TBytes of data. The reads and writes must exceed the disk system cache by at least a factor of 10 in size. The disk system shall have a normal load of existing data. |
| 4.5 The minimum sustained aggregate transfer rate from the disk storage system (classified or unclassified) to its service and I/O node partition while it is de-coupled from the compute partition shall be at least 50 GBytes/s. |
| 4.6 The system must include a minimum of ten 1.0 Gigabit Ethernet network connections for each service and I/O partition, classified and unclassified, to be used for user access to the system. |
| 4.7 The system shall include sufficient 10.0 Gigabit Ethernet network connections to meet the aggregate sustained network I/O bandwidth requirement of 25 GBytes/s. As stated in Table 2, the sustained bandwidth of 25 GBytes/s shall be achieved using the TCP/IP protocol. |

# Architecture and Performance Requirements for System Management and Reliability, Availability, and Serviceability (RAS)

The architecture and performance requirements for system management and RAS are presented in Table 5. These requirements are extremely important for the successful operation of the Red

# Table 5:
# Architecture and Performance Requirements for System Management and Reliability, Availability, and Serviceability

| Description |
| --- |
| 5.1 There shall be two system management workstations, one for managing the classified section and one for managing the unclassified section of the machine. There shall also be a backup for each of these such that failure of a management workstation will not result in an interrupt of either the classified or unclassified sections of the machine. |
| 5.2 System management shall provide a single system image of both sections, classified and unclassified, of the machine. All machine resources (for example processors, I/O, and disk storage) within a section, classified and unclassified, shall be managed from these system workstations. |
| 5.3 A complete Red/Black reconfiguration of the system shall be accomplished in less than one hour. Within the time-limit of one hour the section of the machine to be moved from classified to unclassified or unclassified to classified shall be shutdown, all cable changes shall be made, the memory in the section being moved shall be scrubbed, and the full system (classified and unclassified) shall be brought back up including rebooting where necessary and be ready for user applications to be loaded. Sections of the machine not being switched shall remain operational during the switching process. Jobs running on the unswitched sections shall continue uninterrupted. (Meeting this requirement may require a full system reference clock to maintain synchronization of all of the compute nodes across the whole Red Storm system.) |
| 5.4 There shall be separate, fully independent system management/RAS networks for both sections, classified and unclassified, of the machine, one of which shall include the switchable center section. These networks shall provide system management and RAS access and monitoring for all significant components in the system including the service and I/O node partitions and the disk storage system. |
| 5.5 There shall be dedicated RAS nodes which have their own processors and connection to the system management/RAS network to monitor and manage compute nodes. Each RAS node shall be responsible for a maximum of 32 compute processors. Failure of a RAS node shall not cause a system interrupt. |
| 5.6 In general, the system shall be designed to prevent single-point failures of either hardware or system software that can cause an interrupt of either the classified or unclassified systems. Single-point failures that are possible but very unlikely are acceptable, however, there shall not be any single-point failures that can cause a system interrupt for high failure rate components such as power supplies, processors, compute nodes, 3-D mesh primary communications network, or disks. It is acceptable for the application executing on a failed processor or node to fail but when this happens applications executing on other parts of the system shall not fail. |

# Table 5:
## Architecture and Performance Requirements for System Management and Reliability, Availability, and Serviceability - Continued

| Description |
| --- |

5.7 Each section, classified and unclassified, of the machine shall have a primary and a backup boot node. The primary and backup boot nodes shall be cross linked to the respective, classified and unclassified, boot raids. There shall be an automatic fail over mechanism to prevent a system interrupt due to the loss of a boot node.

5.8 Scalable RAS software shall provide real-time monitoring and tracking of all significant hardware components in the system. These components shall include power supplies, processors, memory, interconnect, communication network interfaces and switches or routers, interconnect cables, service and I/O nodes, disks, disk controllers, and fans. All error conditions, recoverable and non-recoverable, shall be monitored and tracked by component. This RAS software shall include a scalable, graphical user interface running on the system management workstation which provides for display of this data.

5.9 Disk storage controllers and service and I/O nodes shall be redundant and have an automatic fail-over capability. Card cage power supplies shall have a minimum of an N + 1 configuration for each card cage in the compute node partition such that the loss of an individual power supply shall not cause an interrupt of any compute node.

5.10 The Red Storm full system Bit Error Rate (BER) for non-recovered errors in the 3-D mesh primary communications interconnect shall be less than 1 bit in $10^{21}$. This error rate applies to errors that are not automatically corrected through ECC or CRC checks with automatic resends. Any loss in bandwidth associated with the resends would reduce the sustained interconnect bandwidth and must be accounted for in claimed sustained bandwidth for the Red Storm interconnect. Depending on the type of cables used for the Red Storm system (optical versus copper), ECC and/or CRC across each link in the interconnect will probably be required to achieve the required BER. End-to-end CRC checking with automatic resend will almost certainly be needed to reach the required BER.

5.11 A full system reboot of the classified or unclassified sections from a clean shutdown and without a disk system fsck shall take less than 15 minutes.

5.12 Hot swapping of failed Field Replaceable Units (FRUs) is highly desirable. Where hot swapping is provided it shall be possible without power cycling the full system. The maximum number of components (such as nodes, disks, and power supplies) contained in or on one FRU shall be less than 1% of the components of that type in the compute partition of the Red Storm system.

# Table 5:
## Architecture and Performance Requirements for System Management and Reliability, Availability, and Serviceability - Continued

| Description |
| --- |
| 5.13 Mean Time Between Interrupt (MTBI) for the full system shall be greater than 50 hours for continuous operation of the full system on a single application code. This means that the full system must be able to run continuously on an application that is using the full system for 50 hours without any hardware component failures or system software failures that cause an interrupt or failure of the application code. (In this context "full system" means the maximum configuration for the compute partition plus one service and I/O partition.)<br><br>5.14 MTBI for the full system, as determined by the need to reboot the system, shall be greater than 100 hours of continuous operation. This means that the system will be continuously operational for 100 hours with at least 99% of system resources available and all disk storage accessible.<br><br>5.15 Diagnostics shall be provided that will, at a minimum, isolate a failure to a single FRU. This diagnostic information must be accessible to operators through external network connection to the System Management/RAS workstation for the classified and unclassified sections respectively.<br><br>5.16 FRU (or node) failures shall be able to be determined, isolated, and routed around without system shutdown. The operators shall be able to reconfigure the system to allow for continued operation without use of the failed FRU (or node). The capability shall be provided to perform this function from a remote network workstation.<br><br>5.17 There shall be a scalable diagnostic code suite that checks processors, RAM memory, network functionality including NIC and Switch chips and cables, I/O interfaces, and disk controllers for the full system in less than 30 minutes.<br><br>5.18 All bit errors (from memory, interconnect, and disks), over temperature conditions, voltage irregularities, fan speed fluctuations, and disk speed variations shall be logged in the RAS system. All bit errors shall be logged for recoverable and non-recoverable errors. |

Storm system. Computers of this size are difficult to manage, maintain, and operate. Without an integrated, full system design for system management and RAS it would be extremely difficult to meet the MTBI requirements in Table 5. Also, without the ability to track and monitor the behavior of the system hardware components it is nearly impossible to figure out which parts have failed or which parts are likely to fail and should be replaced in the system before they fail.

# System Software Requirements

The system software requirements are given in Table 6. They are based on the need to have a

## Table 6:
## System Software Requirements

| Description |
| --- |
| 6.1 The compute node partition of the system shall have a highly scalable and efficient Light Weight Kernel (LWK) operating system. (The Red Storm LWK shall be derived from Sandia's Puma/Cougar operating system.) The compute node operating system shall not provide support for demand paging, time-sharing, sockets, or other non-scalable UNIX functionality. The LWK operating system shall provide for virtual memory addressing per compute node and multiple processes per compute node. |
| 6.2 The service and I/O nodes shall have a full UNIX or UNIX derivative operating system. The performance / availability of "UNIX" services must scale with the floating-point performance of the system. Where these services are provided by more than one processor, they shall appear to applications as if they are provided as a single scalable service. |
| 6.3 The RAS nodes shall have a real-time or real-time like operating system that supports the RAS network and system management functions. The system management/RAS workstations shall have a full UNIX or UNIX derivative operating system. |
| 6.4 System software shall provide I/O access for applications running on any node of the compute node partition or service and I/O node partition of the system to all secondary storage resources connected to the system. System software will also provide for scalable parallel access to I/O storage resources from applications running on any compute node or service and I/O node of the system. |
| 6.5 Message Passing Interface (MPI-2 or later) standard message passing software library shall be provided. If MPI is an interface to a native message passing library, then the native message passing library must also be provided. The latest generally accepted version of MPI shall be provided. MPI shall be provided for intra-node as well as inter-node communication for nodes that have more than one processor. The MPI provided shall support message passing between heterogeneous executables within a single job. |
| 6.6 System software shall include the latest standard version of NFS. At a minimum NFS shall provide support for serial and parallel file sizes of at least 100 TBytes. |
| 6.7 File system software for both serial and parallel files shall provide for file sizes of at least 100 TBytes. |

# Table 6:
# System Software Requirements - Continued

| Description |
| --- |
| 6.8 Accounting of user resource consumption shall be provided by the system. This accounting shall include the number of processors used, disk storage usage, file I/O, and the application run-time. A log containing this accounting data shall be maintained and accumulated over time. The accounting software shall provide monthly accumulation of overall system availability and utilization. |
| 6.9 The system software shall include optimized run-time libraries for BLAS levels 1,2, and 3. In addition, an optimized set of libraries for serial and parallel dense and sparse matrix operations, math functions including complex arithmetic, and I/O shall be included with the system. |
| 6.10 Parallel file system software shall be provided as part of the service and I/O node system software that provides for striping of files across RAIDs. This software shall be able to sustain a minimum of 50 GBytes/s (reading and writing) for a single parallel file which is being written to by at least 25% of the compute nodes in the full machine. This software shall also be able to sustain 50 GBytes/s (reading and writing) for the case in which a separate file is being written by each compute node. It shall also be able to sustain 50 GBytes/s for situations in which there is an application mix of these conditions. |
| 6.11 The system software shall provide protection to prevent user applications from corrupting systems code and other users' codes. |
| 6.12 The system software shall provide for space sharing of the compute node partition and space and time-sharing of the service and I/O node partition. The smallest size for user allocation of the machine shall be at the level of a single compute processor. |
| 6.13 The compute node and service and I/O node operating systems shall provide for multitasking on processors. |
| 6.14 The full system must be able to be booted in less than 15 minutes after a clean shutdown. |
| 6.15 A **heterogeneous** application (5 MBytes or less per load module and 3 different executables) shall be able to be fully loaded and running at full application speed on the full machine in less than 1 minute. Job loading shall be from a single file system in the service and I/O partition. |
| 6.16 Current ANSI standard, optimizing FORTRAN F90, C, and C++ compilers and run-time libraries shall be provided. These compilers must produce executable code for the Red Storm system while compiling on at least one of the following workstation platforms: Sun SPARC, DEC Alpha, HP PA-RISC, SGI MIPS, IBM Power, or AMD Opteron (Sledgehammer). The cross compiling workstation shall be running a UNIX or UNIX derivative operating system. |

# Table 6:
# System Software Requirements - Continued

| Description |
| --- |
| 6.17 The system software shall include the Totalview interactive, parallel symbolic debugger. Totalview as implemented on Red Storm shall be compatible with each of the above compilers and languages. Totalview shall be able to work with application codes that have sections written in one or more of these languages. Totalview shall also be able to work with a heterogeneous application (multiple executables in a single job). |
| 6.18 The system software shall include a scalable performance monitoring tool. This tool shall have a graphical user interface that runs in the service and I/O node partition and interfaces to monitoring software running in the compute node partition. |
| 6.19 The system software shall support interactive and batch partitions. For the batch partition PBS or an equivalent batch execution software package shall be provided. This package shall provide for scheduling of jobs based on priority, job size, and job run time. Similarly for the interactive partition the system software shall provide for scheduling of jobs based on priority, job size, and job execution time. |
| 6.20 Libraries intended to be used on any node that has more than one processor shall be multithreaded to the level of at least the number of processors on that node. All multi-threaded libraries shall provide protection from corruption through multiple use of the library by separate processes (thread safety). |
| 6.21 A high performance TCP/IP implementation shall be included as part of the operating system software for the service and I/O nodes. The TCP/IP implementation shall be able to sustain a bandwidth of at least 500 MBytes/s per external network interface. |
| 6.22 Support for DFS or an equivalent shall be provided. This requirement shall be met by running the DFS or equivalent software on the service and I/O partition. This is not a requirement for DFS or equivalent on the compute partition. |
| 6.23 Support for HPSS shall be provided. HPSS support may be provided through the use of one or more separate servers. |
| 6.24 Parallel FTP support shall be provided. |
| 6.25 Documentation for all software provided by the contractor with the Red Storm system shall be provided to Sandia. |
| 6.26 Source code shall be provided to Sandia for all system software provided by the contractor for the Red Storm computer system. This requirement does not include compilers or third-party software provided with the machine where the contractor does not have the right to provide source code. |

system that is highly scalable and yet provides for the functionality needed by our applications. Functionality that inhibits good scalability has been purposefully rejected for the compute node operating system while full functionality has been specified for the service and I/O node operating system. Requirements for tools such as compilers, debuggers, and performance monitors and libraries are aimed at providing needed programming environment functionality for application code development and use.

The administrative software requirements are aimed at needed functionality only. A simple batch queuing system such as PBS was chosen because it provides the needed functionality. System level checkpointing is neither required nor desired because it is very inefficient. The system cannot know what data in an application is needed to restart the application and as a result it must write out the entire application code image. This is typically several times as much data as is needed. It wastes systems resources for I/O and for disk storage. Once an application is scheduled into the machine it is assumed to run to completion or until it reaches its time limit.

## Secure Computing Requirements

Sandia has special requirements related to security. These requirements are given in Table 7. Several of these requirements are partially included or impacted by some of the architectural

### Table 7:
### Secure Computing Requirements

| Description |
| --- |
| 7.1 The Red Storm system must support the needs of both classified and unclassified computing. This requires that there be a minimum of two independent file storage systems, one for classified and one for unclassified. Each file system shall have a minimum of 120 TBytes of storage capacity. |
| 7.2 The center section of the compute partition shall be able to be switched from the unclassified to classified mode in 1.0 hour maximum and from the classified to the unclassified mode in 1.0 hour maximum. For the case when the whole compute partition is to be classified or unclassified the reconfiguration shall be accomplished in a maximum of 2 hours. |
| 7.3 There shall be a minimum separation of 1.0 feet between the classified and unclassified sections of the system. This separation requirement applies to all electronic components in the system including but not limited to interconnect cables, disk storage systems, and service and I/O partitions. |
| 7.4 The system must provide Kerberos style ticket forwarding, post dated tickets, and ticket renewal. |

# Table 7:
## Secure Computing Requirements - Continued

| Description |
|---|
| 7.5 A special warranty for classified disks shall be included. Because of the need to protect classified information, classified disks that have failed cannot be returned but must be destroyed at Sandia. |

requirements specified for the computer system. For example, the architectural requirements for Red/ Black switching and the need to support both classified and unclassified computing in a single computer system shows up in the separation requirements given in 7.3 above and the reconfiguration time requirements in 7.2 above.

## Maintenance Support Requirements

Maintenance support will be provided for the system on a seven day a week 24 hours per day basis from the time first significant hardware is delivered to Sandia through the end of the contract. (For the purpose of on-site support the first significant hardware delivery to Sandia will be the installation at and release to Sandia for user application code use of the first 1/4 of the Red Storm computer system.) On normal Sandia business days support personnel shall be on-site from 7:00am till 7:00pm. Outside of these hours and on holidays and weekends support personnel shall be available on call to perform maintenance with a two hour response time. A supply of spare parts will be maintained on-site at Sandia which is sufficient to cover any and all reasonably expected component failures. If returning the Red Storm system to normal operating status requires parts to be shipped to Sandia, parts are to be shipped same day air or, if that is impossible, then next day air.

Regular system software updates shall be provided as part of the system maintenance throughout the life of the contract. These updates shall include operating systems, compilers, debuggers, performance monitors, RAS software, math and I/O libraries, communication libraries, and any other software provided with the system or as part of this contract.

# Task II:Program Management and Support for Task I

The contractor shall execute a Program Management Plan (PMP) consistent with the scope and overall objectives of this program as it applies to Task I. (Table 8 lists minimum project plan parameters.) The PMP shall include an initial Program of Record (POR) for the first year and milestones for Tasks I and II. Each quarter, the current government fiscal year POR will be reviewed and supplemented with a detailed plan for the next two quarters. At the end of the first year, the second year POR will be reviewed for approval by Sandia. At the end of the second year, the third year POR will be reviewed for approval by Sandia. This process will continue through the life of the contract.

## Table 8:
## Red Storm Project Management Plan Parameters

| Description |
| --- |
| 8.1 The PMP shall discuss how the Red Storm development will support the Stockpile Stewardship Computing goals and not be an end in itself. |
| 8.2 The PMP shall discuss how the system capability will be practical to operate, use, and maintain. |
| 8.3 The PMP shall include detailed schedules and development plans for the design and development of the Red Storm computing system including, but not limited to, all major system components to be built or purchased. |
| 8.4 The PMP shall include test plans for building and integrating the Red Storm computer system. |
| 8.5 The PMP shall include plans for installing and integrating the Red Storm computer system at Sandia. |
| 8.6 The PMP shall include a discussion of risks, both technical and schedule, associated with the project and a strategy for mitigating those risks. |
| 8.7 The PMP shall address the life-cycle of the system and the contractor support that will maintain the capability on-site and on-line for at least 4 years after final installation. This includes preventive maintenance, test fixtures on-site at Sandia, procedures for "burning" in spare parts, and other procedures for ensuring maximum availability of Red Storm to the ASCI user community. |

On signing of a contract with Sandia, the contractor shall provide to Sandia the names of the project manager and his team for managing the Red Storm project. Also, at this time the contractor shall provide to Sandia a contract management plan that includes a problem identification and conflict resolution processes.

During the development phase of the contract Sandia and the contractor shall hold weekly conference calls to review progress and to discuss issues or problems that might arise during the development. In addition, during the development phase, Sandia will appoint a project leader for this project that will participate in regular meetings at the contractor's site with the contractors technical project team.

In addition to the above meetings and conference calls, the contractor shall hold monthly project review meetings with Sandia during the development phase. Also during the development phase

of the contract there shall be formal quarterly project review meetings between Sandia and the contractor. After the system has been installed at Sandia and has moved into a stable operating mode, the meeting and conference call schedule will be relaxed.

The contractor shall identify significant risks associated with meeting contract schedules or performance requirements and develop a plan for risk mitigation. The risk mitigation plan shall be updated regularly, at least quarterly. The risk mitigation plan shall be incorporated in the yearly plans. The contractor shall promptly notify Sandia of any issues that might affect the viability or schedule for the project. Sandia and the contractor will work together to manage risks and to achieve the overall project goals.

The contractor shall provide 3 on-site computational scientists for support of the system for the life of the contract starting at the time of first delivery of significant hardware (delivery of the first 1/4 of the system) to Sandia under the contract. The computational scientists shall provide support to the application code developers and other users of the Red Storm system and other system support. The contractor will also supply additional on-site personnel for maintaining and operating the Red Storm system and any other hardware delivered to Sandia as a part of this contract. On site personnel will be directly accountable to Sandia. All on site support personnel will be subject to review and approval as to their qualifications by Sandia and shall be able to obtain a Q clearance.

DISTRIBUTION:

| | | |
|---|---|---|
| 1 | MS-0139 | Mike Vahle, 9900 |
| 1 | MS-0139 | Robert K. Thomas, 9904 |
| 1 | MS-0321 | William J. Camp, 9200 |
| 1 | MS-0310 | Robert W. Leland, 9220 |
| 1 | MS-0310 | Mark D. Rintoul, 9212 |
| 1 | MS-0316 | John Aidun, 9235 |
| 1 | MS-0316 | Sudip Dosanjh, 9233 |
| 1 | MS-0318 | Jennifer Nelson |
| 1 | MS-0318 | Paul Yarrington, 9230 |
| 1 | MS-0819 | Randy Summers, 9231 |
| 1 | MS-0820 | Pat Chavez, 9232 |
| 1 | MS-0822 | Phil Heerman, 9227 |
| 1 | MS-0847 | Ted Blacker, 9226 |
| 1 | MS-0847 | Scott Mitchell, 9211 |
| 1 | MS-1109 | Paul Iwanchuk, 9220 |
| 1 | MS-1110 | Doug Doerfler, 9224 |
| 1 | MS-1110 | Neil Pundit, 9223 |
| 1 | MS-1110 | David Womble, 9214 |
| 1 | MS-1111 | Bruce Hendrickson, 9215 |
| 1 | MS-1109 | James L. Tomkins, 9220 |
| 1 | MS-9018 | Central Technical Files, 8945-1 |
| 2 | MS-0899 | Technical Library, 9616 |