

Dark Storm: Further Adventures in XT architecture flexibility

**John P. Noe, Robert A. Ballance, Geoffrey McGirt,
Jeffry Ogden, Sandia National Laboratories**

ABSTRACT: *The Cray/Sandia XT3 architecture instantiated in the Red Storm computer system at Sandia National Laboratories provides many opportunities for upgrade, modification and customization. The original 40TF platform has seen several upgrades over the years resulting in its current 284TF peak configuration. The system has also received updated disk subsystems reaching over 2 PBytes of storage split between the two original network heads. Recently Sandia was tasked with supporting efforts within DOE/NNSA aimed at extending the reach of traditional supercomputing into non-traditional areas of interest to National Security organizations. To support this new initiative, Sandia and Cray combined concepts to create, deploy and demonstrate an external network supported Lustre file system routed to Red Storm via multiple 10 GE connections. Access to Catamount compute nodes is routed through Service Nodes running the LNET router protocol. This concept was deployed and initial demonstrations indicated file system throughput of over 12 GB/second and support for over 12000 nodes. Utilizing an innovative Woven network switch which provides load balancing message handling, the concept can be extended to support multiple customer specific file systems all with equal access to the entire Red Storm system. Such separate file systems would connect serially to Red Storm, but the data would remain available to customers whether connected to Red Storm or separated. This talk depicts the testing, validation and debugging efforts which resulted in the successful deployment of this National Security resource, and some of the issues associated with debugging in a high security environment.*

KEYWORDS: Red Storm, XT, Lustre, National Security, Catamount, LNET, CLE

1. Introduction

Red Storm

The Red Storm computer platform, conceived by Sandia National Laboratories and built by Cray, Inc to address critical needs in the National Nuclear Security Administration Office of Advanced Simulation and Computing was the genesis of the Cray commercial XT series of computer systems.

Red Storm remains unique among the nearly eighty installed XT3/4/5 systems with its Red/Black switch capability, and a two headed configuration. Sandia continues to use Catamount, our Light Weight Kernel operating system on the 12,980 compute nodes, making the system not quite like other deployed platforms. The LWK helped make Red Storm extremely scalable and served as a benchmark for CLE noise reduction efforts within Cray.

The system exists in two classification regimes: Restricted and Secret. The Restricted network is the main development environment for Sandia researchers and

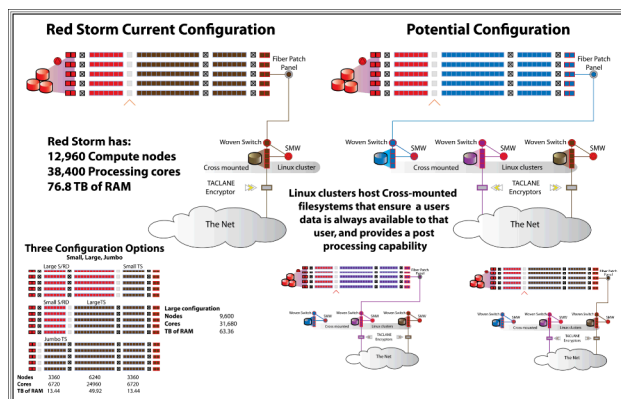
code development teams. The Secret network provides for operations supporting Nuclear Weapon's analysis within the Defense Programs mission space. Red Storm has served the needs of NNSA/ASC programs at Sandia, Los Alamos and Lawrence Livermore National Laboratories since its initial operational capability in March 2005.

Red Storm as originally configured was approximately 40Tflops peak, with 10,368 single processor AMD Opteron nodes in the compute partition. An initial upgrade in 2007 to dual processors, and a fifth row of 31 cabinets (4 I/O and Service cabinets, 27 compute), was followed by another upgrade of quad-core processors to the 65 center section cabinets in 2008. Red Storm is currently equipped with 38,400 processors with a peak of 284Tflops.

Dark Storm

The mission of Sandia National Laboratories is evolving from the basic Nuclear Weapons mission space toward more diverse activities focussed on National Security issues in Energy, Non-proliferation, Homeland Defense and Cyber Security. The traditional focus of our supercomputing systems was to support high performance computational research in science and engineering with an emphasis on nuclear weapons stockpile stewardship. Today, Sandia provides high performance computational, visualisation and data storage resources for researchers and analysts supporting traditional Department of Energy, National Nuclear Security Administration programs, and has recently been expanding support within the Department of Defense, DOE Office of Science and other academic institutions and government agencies. Sandia is located on Kirtland Air Force Base in Albuquerque, New Mexico.

Sandia has always been ready to respond rapidly to requests for assistance in matters of National Security. During a few week period in January and February 2008 Sandia dedicated Red Storm to a high priority mission called Operation Burnt Frost. The mission details were closely guarded and we did not learn the nature of the tasks until several months latter when it was announced that Red Storm had provided the confidence to assure destruction of a disabled satellite which was shot down from its deteriorating orbit by the Navy. This success and Sandia's responsiveness laid the groundwork for an emerging mission space in National Security for Red Storm. Thus, Dark Storm the project was born.



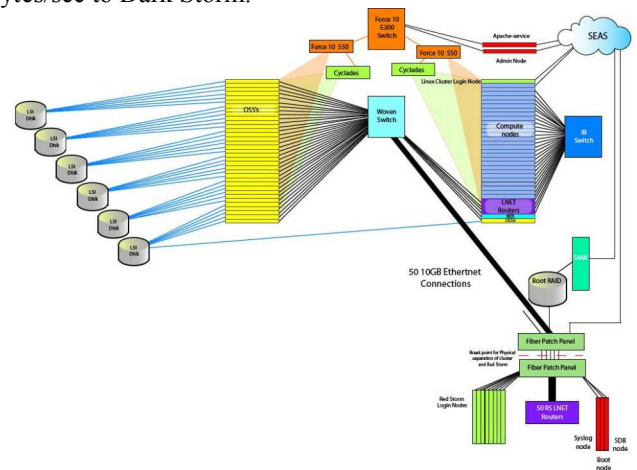
2. Description of systems - hardware

The challenge placed before Sandia was to conceive and demonstrate a configuration for Dark Storm that would permit the essential Nuclear Weapons work to continue on Red Storm, while simultaneously providing confidence to interested sponsors that the full power of

the system could be deployed in a new, high security environment. Merely repurposing the existing restricted network portion of Red Storm would not be sufficiently compelling, and would provide solutions for only one customer. The limitations of an internally hosted file system had to be overcome.

Lustre routers were in use on several Sandia capacity cluster systems, including the 4480 node Thunderbird system. This gave use some confidence that a network attached Lustre file system could be made to work for Dark Storm. Combining this concept with an innovative 10GE high performance network switch built by Woven provided the basic architecture for the prototype demonstration.

An older Linux cluster was repurposed to serve as the Lustre file system host. Disks from the now retired restricted portion of Red Storm were reused. LNET routers were built and installed on the Service Nodes of the system, connected through the Woven switch to the Lustre cluster *Apache*. As many as 50 connections at 10 Gigabits/second are available from Dark Storm. Thirty six 10 GigE connections were made to Apache. This number was derived from the number of Lustre OSSs. The design goal for file system performance was 20 GBytes/sec to Dark Storm.



3. Software

3.1 Operating System

The Catamount light weight operating system is used on Red Storm and Dark Storm. Catamount interfaces with Cray SLES operating systems running in the Service partition. Most operating system functionality dealing with external connections and interactive response for customers is centered in SLES, not Catamount. The non-standard behaviour of Catamount, specifically avoiding interrupt handling from network connections, is in conflict with assumptions made by Lustre. The LibLustre

library runs in user space on Catamount and acts as the interface between the applications on compute nodes and the Lustre file system space. In Red Storm, as in other XT systems, the Lustre File system is running on local disk and is mounted directly across the high speed interconnect network. Red Storm provided separate file systems for each network head, and demonstrated peak performance of 46 Gbytes/sec to applications.

Cray Linux Environment operating system has been competing against Catamount for scalability for some time and appears nearly equivalent for many if not most applications. However, the bulk of codes running on Red Storm are Catamount based and would require extensive effort to validate under CLE. We wanted to avoid putting this burden on the code teams if possible. Thus, a Catamount based solution was preferred.

3.2 File systems

The Lustre file system, Version 1.4 is supported under Catamount pre-2.0 releases and was the current version when our transition began. Although we were aware of the difficulty in supporting an “end of life” version of Lustre, the implications of moving to Lustre 1.6 were beyond our team resources. The client side support offered by LibLustre would have to be integrated and tested with Catamount at full system scale. Our LNET experience on Lustre 1.6 servers and clients in our capacity cluster systems was good but we were aware that release 1.8 was planned and the incremental improvement in support for moving to 1.6 was judged insufficient for the work involved. Thus, we constructed a Lustre 1.4 server system to support the existing Catamount client software which had been demonstrated at full scale on Red Storm, albeit with the internal Lustre.

LNEXT integration and testing proceeded over several months with dedicated support from internal Sandia Lustre experts as well as Cluster File Systems personnel and Cray software support. Eventually, the system demonstrated sufficient stability (although with a few quirks to be addressed) to attempt an application study up to the full extent of the SMALL configuration: 6720 cores on 3360 nodes.

This testing period was conducted in an unclassified environment which made interacting with the support personnel relatively trouble free. Although, sending log files of several 10s of gigabytes presented a challenge. Working with an International team of analysts actually proved viable. Logs gathered during testing in New Mexico were transmitted overnight and analyzed early the next day in Europe, followed by suggestions for change which were ready by western U.S. working hours and the cycle repeated until the system passed testing.

3.3 Compilers, tools, libraries, multi-processing support

The standard tools and libraries extant on Red Storm were pushed to Dark Storm virtually unchanged.

Binary versions of applications were utilized initially as investments in compilation environments had not been made for the Dark Storm configuration. Initially, Sandia applications were the only user codes run. Eventually, a cross compiler workstation was implemented to facilitate introduction of additional codes from collaborators. Ideally, new codes are tested on Red Storm prior to being moved to Dark Storm.

Compilers and debuggers are adequate for the most part yet scalability remains a question and presents unique challenges when operating in a high security environment. Customers are very concerned over matters of data integrity and segregation from other users.

3.4 3rd party applications

Some new codes have been proposed for use on Dark Storm. These include some which have never been ported to Red Storm or the Catamount environment. During the prototyping effort, some codes were ported to assess the difficulty of converting to the Catamount system. Results were generally favorable. As yet, no code has been presented which would force a conversion to CLE due to incompatibility with Catamount.

3.5 Network issues

Our high performance network analysts helped configure and install the Woven switch and analyzed early performance to ensure full bandwidth performance. Sadly, the Woven producer has been bought out and further products are uncertain. Discussions are underway with Arista to obtain and test a similar switch.

3.6 peripherals - disc, tape

Both DDN and LSI disk subsystems are supported on Red Storm. The initial cadre of DDN 8400 controllers has proved problematic from a field failure rate perspective. LSI provided the disks procured under a competitive upgrade purchase. Red Storm can transfer files to an HPSS system for long term storage. Dark Storm has a more limited infrastructure available and thus has no network attached home user files or long term backup via tape systems..

4. Operations

4.1 Resource allocation

Dark Storm is operated in a single user mode, but it could support multiple team members if they are all working on the same project. Need to Know demands are extremely complex in this new operating environment so we are taking our time developing additional operating modes. Full utilization of the system is less valuable than protecting the information being analyzed and ensuring data purity and integrity to our new customers.

Sandia uses the standard MOAB scheduler and Torque allocator under Catamount.

4.2 Data management, backups

File system space management is a peer pressure based ad-hoc managed system. Administrators have not had to intervene with purges at any point in the life of Red Storm. The smaller capacity of Dark Storm is outweighed by the single project nature of the work. There will be incentive for the customers to police their own use to provide adequate space for analysis runs.

5. Applications performance

5.1 Input/Output

Limited today by the size of the Lustre support cluster, we are hopeful that the only limitation on I/O performance is hardware resources, whether implemented under the Lustre/LNET/Catamount model or with CLE and DVS. Panassas is also a viable candidate external file system which will be demonstrated shortly by the Cielo platform, a Los Alamos and Sandia collaboration.

5.2 Execution

Applications codes operate on Dark Storm at exactly the same rate as on Red Storm. Performance has not been impinged by the new configuration.

6. Usability

6.1 Documentation

All Dark Storm user documentation is being transferred into the appropriate security environment but will remain available on the classified environment as well. This model, where lower security levels are utilized in a progression to higher levels has proven beneficial to administration of the system, as well as application development and testing.

6.2 Error handling

Errors? We don't get no stinking errors!

6.3 Vendor support

Sandia wishes to acknowledge the support of Cray and the onsite Cray staff during this experiment in developing a new operating configuration as we seek to secure additional customer sponsors and expand the relevance of high performance computing into new and underserved areas of interest. Dick Dimock, Barry Oliphant, Robert Purdy, Jason Repik and Victor Kuhns have been coddling and cajoling Red Storm since its inception and have seen it through many ups and downs. We could not do this without their professionalism and attention to detail.

7. Observations

7.1 Price of admission

The new customers benefit from the capital investments provided by the Nuclear Weapons Advanced Simulation and Computing program and avoid the initial acquisition cost while they evaluate the usability of these resources. Hopefully, successful mission accomplishment will provide the impetus to recapitalize the Dark Storm system and supply HPC resources into the future.

7.2 Sustained performance measures

Sustained performance is hard to define – it depends on the application and the nature of the work. Dark Storm provides answers to critical questions in hours to days, versus weeks to months from existing computer systems. Often this is the crucial difference in taking action or missing out on mission success.

Conclusion

Both the Red Storm and the Dark Storm environment have demonstrated the usefulness of the XT3 design flexibility and the advantages of the Sandia Red/Black switch concept. The ability to deploy on small systems has been very useful during the transition and the ease with which systems can be reconfigured also enhances the operational use of the systems and the flexibility to offer system time and operational modes which support maintenance and upgrade activity. Sandia is pleased that Dark Storm is poised to lead yet another advancement in pursuit of Exceptional Service in the National Interest.

Acknowledgments

The authors would like to thank Sandia colleagues and the Office of Advanced Simulation and Computing for the opportunity to create and deploy exciting technical solutions for National Security.

About the Authors

John Noe is Manager of Scientific Computing, Center for Networking and Computing, Sandia National Laboratories. He is a long-time CUG member and formerly served on the CUG Board of Directors. He can be reached at P.O. Box 5800, Dept. 9328, MS0807, Albuquerque, New Mexico, USA, 87185-0807. E-Mail: jpnoe@sandia.gov. Robert A. Ballance, PhD, is a Distinguished Member of Technical Staff at Sandia and has been the System Manager for Red Storm since it's inception. Geoffrey McGirt and Jeffry Ogden were the technical leads in developing and deploying the Dark Storm demonstration environment. All may be reached at the same department address above. E-mail: raballa@sandia.gov, gcmcgir@sandia.gov, jbogden@sandia.gov.

Sandia National Laboratories is a multi-program laboratory operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin company, for the U.S. Department of Energy's National Nuclear Security Administration.

NOTICE: This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government, nor any agency thereof, nor any of their employees, nor any of their contractors, subcontractors, or their employees, make any warranty, express or implied, or assume any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represent that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government, any agency thereof, or any of their contractors or subcontractors. The views and opinions expressed herein do not necessarily state or reflect those of the United States Government, any agency thereof, or any of their contractors.