



Deployment and testing of the sustained petascale Blue Waters system



Celso L. Mendes*, Brett Bode, Gregory H. Bauer, Jeremy Enos, Cristina Beldica, William T. Kramer

National Center for Supercomputing Applications, University of Illinois, Urbana, IL 61801, United States

ARTICLE INFO

Article history:

Received 1 October 2014
Received in revised form 20 March 2015
Accepted 23 March 2015
Available online 2 April 2015

Keywords:

Large system deployment
Acceptance testing
Petascale performance

ABSTRACT

Deployment of a large parallel system typically involves several steps of preparation, delivery, installation, testing and acceptance, making such deployments a very complex process. Despite the availability of various petascale systems currently, the steps and lessons from their deployment are rarely described in the literature. This article documents our experiences from the deployment of the sustained petascale Blue Waters system at NCSA. Our presentation is focused on the final deployment steps, where the system was intensively tested and accepted by NCSA. Those experiences and lessons should be useful to guide similarly complex deployments of large systems in the future.

© 2015 Elsevier B.V. All rights reserved.

1. Introduction

Blue Waters is one of the most powerful supercomputers currently available for the open-science community. Sponsored by the US National Science Foundation (NSF) and installed at the National Center for Supercomputing Applications (NCSA) of the University of Illinois at Urbana-Champaign, Blue Waters is also the largest machine ever built by Cray. In addition, it has tremendous amounts of memory and persistent storage. Various application groups are achieving the sustained petascale capability of the system, and there is a huge potential for scientific discoveries in the coming years.

This article contains two contributions that are rare in the literature. First, it reveals several details from the machine deployment process, including methods and procedures that were followed for system assessment and acceptance. Second, it provides first-hand lessons that we learned from that deployment, based on the obstacles that we faced and the solutions adopted. These experiences should become useful to guide similarly complex deployments in the future. In addition, the article serves as an early evaluation of Blue Waters: the presented performance results can be used as a reference by the application groups, as they continue to tune their codes to the Blue Waters architecture.

The remainder of this article is organized as follows. Section 2 briefly describes the Blue Waters architecture, and Section 3 presents the timeline of its deployment. Section 4 shows the infrastructure created by NCSA to support the deployment. The acceptance tests and many of their results are presented in Section 5. Section 6 contains details of the post-acceptance upgrade of the system's acceleration capability, Section 7 lists reliability figures observed during initial operations, and Section 8 briefly outlines the Web-based environment created for user support. Major lessons from the deployment are highlighted in Section 9. Finally, Section 10 concludes our presentation.

2. Blue Waters architecture

Blue Waters has an architecture as depicted in Fig. 1. Its computational component is heterogeneous and contains both XE and XK compute nodes; an XE node contains two 2.3 GHz AMD-Interlagos × 86 processors with 16 integer cores per processor, whereas an XK node has one AMD-Interlagos processor and one NVIDIA-Kepler K20X GPU. The XE nodes have 64 GB of memory, while the XK nodes have 32 GB for the CPU and 6 GB for the GPU. Besides the XE and XK compute nodes, the system also has 784 service nodes, which provide functions such as I/O, external access, boot, and others. The total number of nodes, after completion of the 2013 upgrade described in Section 6, is 27,648, hosted in 12 rows of 24 cabinets. All these nodes are interconnected by a single Gemini-based network, with a 3D torus topology of dimensions 24 × 24 × 24. Each Gemini router connects two nodes to the torus.

* Corresponding author. Tel.: +1 2172449234.

E-mail addresses: cmendes@illinois.edu (C.L. Mendes), brett@illinois.edu (B. Bode), gbauer@illinois.edu (G.H. Bauer), jenos@illinois.edu (J. Enos), beldica@illinois.edu (C. Beldica), wtkramer@illinois.edu (W.T. Kramer).

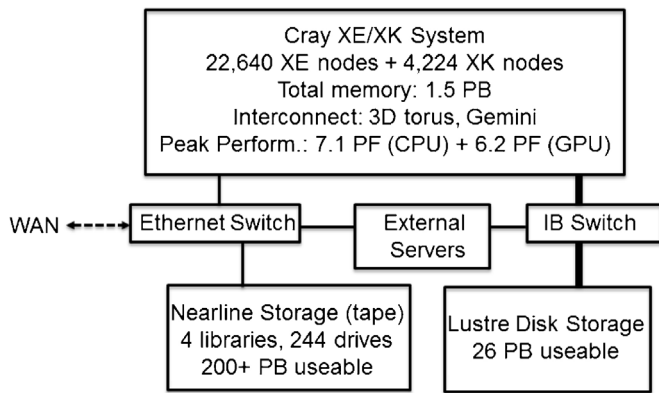


Fig. 1. Blue Waters system architecture.

There are various external servers in Blue Waters, for services like remote login by users, import/export of data, or transfer of data between storage devices. In particular, there are four login servers with the same processors as in the compute nodes; these login servers are aimed at hosting users' editing, compilation, job submission, etc. External connectivity is provided by WAN links with total speed beyond 100 Gbps.

Storage in Blue Waters is available in the form of disk and tape. The disk storage is composed of more than 17,000 hard-disk drives forming three Lustre file systems (`/home`, `/projects` and `/scratch`, with useable capacities of 2 PB, 2 PB and 22 PB, respectively). The aggregate transfer rate to/from the disk storage is more than 1 TBytes/s. Meanwhile, the nearline storage component has four robotic tape libraries with 61 drives per library. The total number of tape slots is beyond 63,000. While the disk storage was supplied by Cray, the tape libraries were provided by SpectraLogic and IBM. These tape libraries comprise the largest High-Performance Storage System (HPSS) installation in the world. The HPSS component also includes 50 servers acting as data movers, and 1.2 PB of disk storage serving as cache to the tapes. More information on the system can be found at the Blue Waters portal (see Section 8).

The site hosting Blue Waters is an advanced datacenter that NCSA built mainly for this purpose. It has several environmental-friendly features, including a LEED Gold certification. Free-cooling can be provided by three cooling towers adjacent to the building; due to the cold weather in Illinois during a good portion of the year, those towers provide nearly all required cooling for several months. For available energy, the building currently has a capacity of 24 MW; actual consumption by Blue Waters has typically been around 10 MW.

3. Deployment timeline

The Blue Waters project started in 2006, in response to the NSF solicitation [1] fostering a sustained petascale system. NCSA was declared the winner of that competition in August 2007. In the following years, NCSA worked with vendors and with various application groups, selected by NSF, to prepare their applications to run at petascale level. As part of the contract signed with Cray, a Statement-Of-Work (SOW) was crafted, containing a sequence of steps for the delivery by Cray of the computational part of Blue Waters. The first of those steps was installation of a Test and Development System (TDS), consisting of one cabinet populated with XE and XK nodes. This machine was installed in December 2011.

Early in 2012, Cray delivered the first 48 cabinets of the final machine, with XE compute nodes. This subset was named the Early Science System (ESS). In the Summer of 2012, Cray delivered additional XE and new XK cabinets (still without GPUs), and integrated them to the ESS parts to form a complete 276-cabinet system. At

the end of that Summer, Cray also delivered the recently-produced NVIDIA-Kepler GPUs and most of the disk storage component, which included new controllers that doubled the bandwidth to disks.

Acceptance tests by NCSA staff started at the end of September 2012, and were the exclusive work on the machine in the following month. In November, use of the system was split between testing and application executions by selected science users. During an NSF review in December 2012, where NCSA submitted to NSF an acceptance report comprising fifteen technical documents describing details of all tests, the Cray system was formally recommended for acceptance. Early in 2013, the other (non-Cray) parts of the system were completed, including the HPSS nearline storage and the full WAN connectivity. Simultaneously, all users approved by NSF were granted access to the system. Blue Waters was officially launched in a ceremony on March 28, 2013, and production usage started a few days later. In the Summer of 2013, 12 additional cabinets with XK nodes were integrated to the system, expanding its torus configuration to $24 \times 24 \times 24$.

4. Infrastructure created for acceptance

In preparation for acceptance of Blue Waters, NCSA staff designed hundreds of tests covering both functionality and performance for all system areas. Those tests encompassed all SOW items and also additional features that NCSA considered important for productive system operation. This section describes the main characteristics of this test design phase and of the infrastructure created for the management of their execution towards system acceptance.

4.1. Test design and preparation

Since Blue Waters has a large number of components with complex interactions among them, the resulting number of tests required for good system coverage is also large. Thus, from the beginning of the project, it became clear that a suitable structure was needed to organize the tests. This structure should be simple to use, but also flexible to accommodate future needs from subsequent phases of the project.

NCSA developed a web-based tool as interface to a database storing the testing structure. This database was intended to serve as a planning and management tool to store numerous characteristics associated with each of the defined tests. The test planning tool was referred by project members as the *Test Matrix* and its GUI is depicted in Fig. 2. In the main panel, on the lower right, each row corresponded to a record in the database. This record was associated with a certain test, identified by a test-ID and a brief test description. In addition, there were several other fields containing information about that test, such as the staff member in charge of the test, the category of the test (e.g. SP = System Performance, ST = Storage, etc.), whether the test required a dedicated system for execution, and other attributes. This view could be customized, and field columns could be added or hidden from the screen as desired. The panel on the lower left allowed filtering to restrict the visualization to records with certain characteristics. At any moment, the panel on the upper right displayed basic counts of records being shown at that moment, according to their status. Finally, the panel on the upper left allowed resetting of the current filters, creation of new records, or generation of specific reports.

In its final configuration, the Test Matrix contained information related to more than 300 tests, as we highlight in Section 5. Many of those tests started to be exercised in the Early Science System. Meanwhile our Test and Development System (TDS) remained, throughout the project duration, our main platform for test development and debugging.

Download English Version:

<https://daneshyari.com/en/article/430362>

Download Persian Version:

<https://daneshyari.com/article/430362>

[Daneshyari.com](https://daneshyari.com)