

Cray Operating Systems Road Map

Charlie Carroll, *Cray Inc.*

ABSTRACT: *This paper discusses the Cray Operating Systems road map. The Koshi (2012), Nile (2012), Denali (2012), Ohio (2013) and Olympic (2013) releases are discussed. In addition, the rationale for coming changes is discussed.*

KEYWORDS: Operating systems, releases

1. Introduction

The Cray Software Operating Systems and I/O (OSIO) group provides key infrastructure and service components of the software stack.

These components include:

- Compute node kernels
 - XT CNL (Compute Node Linux)
 - NVIDIA GPU drivers
- Service node kernel
 - Supports all compute node types
- File systems
 - Lustre
- Networking
 - GNI and DMAPP
 - IBGNI (IB verbs -> Gemini)
 - TCP/IP
- Operating system services
 - Core specialization
 - Dynamic shared libraries
 - Cluster Compatibility Mode
 - DVS (Data Virtualization Service)
- System management
 - Cray Management Services
 - Node Health Checker
 - ALPS (Application Placement Scheduler) Level
 - Command interface
- Third-Party Extensions
 - GPFS
 - Panasas PanFS
 - Batch schedulers

This paper discusses the main themes to be emphasized in upcoming OSIO software releases, followed by specific features to be delivered in these releases.

2. Release Themes

Upcoming OSIO releases will emphasize certain broad themes. Before getting into specifics, we will take a look at the big picture.

2.1. System stability

Stability and robustness are important to any customer system. They are especially important in large supercomputers with millions of separate components.

Cray has invested substantially over the past three years in defect reduction. By our internal measures, we've made substantial progress in reducing customer bugs. By the most important measure—system availability—we've made great progress. System availability is well above 99%.

2.2. Performance

Cray's Compute Node Linux (CNL) implementation performs extremely well, largely because we have limited which services and features run on compute nodes.

In addition, current and upcoming work will focus on GPU performance improvements. One example is host-initiated GPU-to-GPU transfers across the Cray high-speed network.

2.3. Hardware Support

Part of OSIO's mission is to support new Cray hardware as it becomes available. Koshi will add support for NVIDIA's Kepler accelerators in Cray XK systems. Nile and Denali will support Cray's Cascade infrastructure, as will Ohio and Olympic.

2.4. File Systems

Cray supports a variety of I/O models. Much of our installed base uses direct-attached Lustre file systems. Cray Software also ships and supports external Lustre file systems. These integrate white-box servers, Lustre server software and storage. Cray supercomputers connect to these external Lustre servers with Infiniband.

Cray now offers its Sonexion 1300 Lustre appliance.

Cray supports other file systems through DVS (Data Virtualization Service). DVS projects the file system from Cray service nodes to the Cray compute nodes. The service nodes, in addition to serving DVS, are also clients to the remote, projected file system. This enables applications running on the compute nodes to access, through DVS, the external file system. To date, Cray has used DVS to interface with Panasas, GPFS and NFS.

3. Upcoming Releases

Cray will release CLE Koshi in December 2012 for Cray XE/XK systems. CLE Nile will be generally available in March 2013 on Cascade systems. Early releases will support Cascade shipments in 2H12. SMW Denali will support both Koshi and Nile. CLE Ohio and SMW Olympic will release in 2H13. Cray also plans update releases, approximately quarterly, for bug fixes and smaller features.

3.1. CLE Koshi

CLE Koshi will include several new kernel features: Compute Unit Affinity, Compute Node Clean-Up and faster warm-boots.

To explain Compute Unit Affinity (CUA), we need to start with some terminology. Cray will use the term "compute unit" to refer to a group of CPUs which share processor resources like Intel Sandy Bridge and AMD Interlagos processors do. Initial Cray compute nodes with Intel Sandy Bridge processors will have two NUMA nodes, each with eight compute units, where each compute unit consists of two CPUs, for a total of 32 CPUs. A

compute node with AMD Interlagos processors has four NUMA nodes, each with four compute units, where each compute unit consists of two CPUs, for a total of 32 CPUs.

CUA allows users to control how their applications are applied to the compute resources on a node. Specific updates for this feature include:

- modifying the system database to include values for the number of compute units per compute node and the number of CPUs per compute unit;
- updating the `xtprocdadmin` command to display compute unit data for compute nodes;
- updating the `cselect` command to let users select compute nodes based on compute unit attributes;
- modifying `apstat` to display the number of compute units on a compute node;
- modifying ALPS to be compute unit aware;
- adding a '-j' option to `aprun` to allow the user to specify how many CPUs per compute unit should be reserved for the application;
- updating the `libjob` and kernel affinity code to adhere to the `aprun` '-j' option when binding processes to CPUs;
- removing and/or clarifying the term "core" in Cray documentation.

Compute Node Clean-Up (CNCU) will remove, at the end of a reservation, the temporary files and data structures stored in memory on compute nodes. These include files in `/tmp`, `/var`, `/var/lib/hugetlbfs`, sysV IPC (`shm`, `msg`, `ipc`) data structures, and Posix IPC (`shm`, `msg`, `ipc`) data structures. Note that the kernel removes these at the end of the reservation, not the application. That is, these structures can be used to keep data across application invocations.

Faster warm-boots will reduce the time needed to warm-boot a compute node to about one minute.

CLE Koshi will include a Lustre 1.8.7 client, a Lustre 2.2 client, and a Lustre 1.8.7 server. The Lustre 1.8.7 client and server support direct-attached storage. The Lustre 1.8.7 and 2.2 clients support external servers.

CLE Koshi will support Lightweight Log Management (LLM). LLM will provide a standard logging approach and infrastructure based on `rsyslog` that is simple, lightweight and flexible that allows sites to filter and distribute log information to their specifications and leverage third-party tools to mine and report upon that data.

CCM (Cluster Compatibility Mode) and IAA (ISV Application Acceleration) will be improved in CLE Koshi. Many applications will run faster with fewer issues based on these fixes and improvements.

NVIDIA Kepler support will be part of CLE Koshi. In addition, soft resetting of the GPU will be enabled. This can be used by the control subsystem to quickly bring GPUs back into service. In addition, GPU memory can be scrubbed between application runs.

Two application resiliency features will be included in CLE Koshi.

Application Relaunch gives users the option to relaunch their application in place in the event of a node failure. Instead of tearing down the N-node application and exiting, ALPS can relaunch the application, either on N nodes (the user must include one or more spares in the reservation) or on the N-1 nodes. The application needs to manage its own checkpointing and restarting. A major benefit is that the job can run to completion without waiting again to get to the top of the job queue.

ALPS Reconnect allows an application to survive a node failure by rebuilding the application communication tree in place. After rebuilding its tree, ALPS passes information to PMI which rebuilds its tree and then passes information to the programming model. It's up to the programming model to take appropriate action, presumably to divide the work among the remaining processors. We expect that it will be some time before MPI is able to handle down nodes. CHARM++ applications may be the first to survive node drops.

3.2. CLE Nile

Nile, and its companion SMW Denali release, will support Cray's Cascade systems. These include Intel processors (initially Sandy Bridge and Ivy Bridge), Aries (Cray's high-speed network), and a new cabinet infrastructure.

Nile will include all the features described in the Koshi section of this paper. An exception is that only Lustre 2.x and beyond will be supported with external Lustre file servers in Nile.

Five power management features will be offered in Nile: enhanced power monitoring, job power profiling, static system power capping, p-state control at job launch, and idle node power conservation. One, HSS will be enhanced to report power data. Two, power data will be able to be connected with job data. This will allow administrators to tie power usage to a particular job. Three, administrators will be able to set a static power cap

for the entire system. Four, users will be able to set their job's p-state at launch time. Five, idle nodes will be able to set to a lower power state.

3.3. CLE Ohio and Olympic

CLE Ohio and Olympic will release in mid-2013. The feature lists for these releases are still being formulated. Cray expects to support a Cascade accelerator in late 2013.

4. Conclusion

This paper has presented specific features which will be coming in 2012 and 2013 releases of Cray's operating system. In addition, we have discussed the themes and thought processes behind our plans.

Acknowledgments

The author would like to thank his colleagues and development team at Cray. Their commitment to producing the world's best supercomputers makes it a pleasure to come to work every day, as well as making this paper possible.

About the Author

Charlie Carroll is Director, OS and I/O with Cray Inc. If you have comments on our road map, he would love to hear from you at charliec@cray.com.