

The Tri-Lab Operating System Stack on Cray EX Supercomputers

Cray User Group 2023

Trent D'Hooge and Adam Bertsch
Livermore Computing

May 7-11, 2023
Helsinki, FI



Let's Talk About TOSS

- What is TOSS?
- Why TOSS?
- Current Status of TOSS on Cray EX
- Challenges
- Future Work



TOSS is RedHat with Sprinkles on top

The Sprinkles



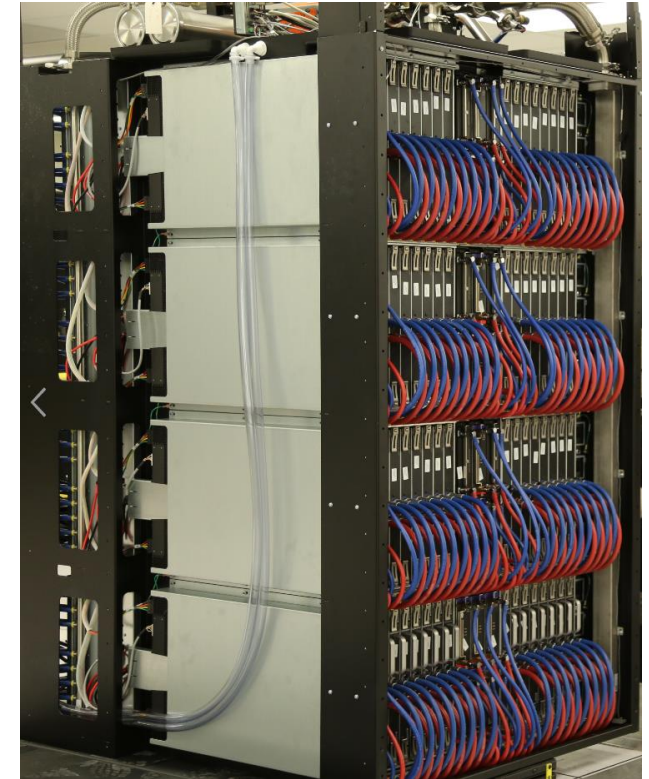
- Release every month
 - Every version of every package is tracked for every release
 - Base OS like RHEL Server and the AppStream “like” stuff
- ~150 Sprinkles:
 - Cassini support
 - Beegfs, GPFS, 3 lustres versions 2.15 with Cassini support
 - Nvidia / AMD GPU support with supporting software
 - SLURM, Flux
 - Cluster management tooling; power, console, parallel shells....
 - Vendor specific tools
 - Dell, SMC, firmware, Proprietary Mellanox OpenSM, Intel
 - Tooling for diags and firmware updates
 - Compilers, MPIs
 - Kernel patches

Why TOSS?

- Simple
 - TOSS philosophy keeps things as simple as they can be in order to improve reliability.
- Mature
 - Over 20 years of incremental development and use on production Top 500 HPC systems.
- Community Oriented
 - Upstream whenever possible. Standard components wherever they exist. Open source any additions where possible.
- Secure
 - TOSS provides monthly security patch releases.
 - TOSS follows a published STIG security baseline.
- Software Supply Chain
 - Every package in TOSS is signed providing a 100% trackable pedigree.

Current Status of TOSS on Cray EX

- Running on four Cray EX systems at LLNL
 - Three production early access systems
 - One testbed
- Full suite of management features
 - Diskless and Diskfull provisioning
 - Power and console management
 - Firmware update
 - Diagnostics
 - Rolling updates
- Shipped on El Capitan infrastructure systems
 - Pre-installed at HPE manufacturing by LLNL and HPE engineers
 - Used to perform pre- and post- ship diagnostics for El Capitan infrastructure
- Base OS for El Capitan ClusterStor filesystem



Challenging Interface Points With Cray EX Software

- Firmware
 - Production level HFP-firmware release schedule
 - To keep consistent update process in place, pre-production firmware is repackaged to match HFP
- Kernel Drivers
 - Delay in compass github repos to match Slingshot release
 - Kernel related software, cxi driver, kabric...
- Cray PE / Modules
 - Early issues have been resolved, packages are production quality
- Monitoring
 - Cray was ahead of the dmtf standards
 - Cray uses OEM endpoints to capture lots of data with one request
 - LDMS support libfabric support?
 - Cassini is not a provider upstream

Future Work

- K8s
 - Extending K8s support from Rabbit hardware to a service on compute nodes
- Using upstream xpmem
 - HPE has opened tickets with upstream
 - Need to find people time to work on issues
 - HPE looking at alternatives
- Moving to RHEL 9
 - TOSS 5.0 work to begin in May 2023 based on RHEL 9.2
 - Likely move to RHEL 9 based TOSS in 9.5 or 9.6 timeframe



**Lawrence Livermore
National Laboratory**