## PDSF Site Report HEPiX Spring 2016





#### James Botts Computational Systems Group

April 18, 2016







- Now Located Lawrence Berkeley National Lab
- NERSC is the primary computing facility for the US DOE Office of Science
- Division of LBNL
- Over 6000 users
- over 400 projects
- 40<sup>th</sup> Anniversary in 2014







# Systems at NERSC – All run HEP workloads



NERSC-8 Phase 1 Cray XC40 1630 Nodes 52160 cores (Haswell) 30 PB Lustre Burst Buffer





NERSC-7 Cray XC30 5200 Nodes 124800 cores 2.4 PFlops Theoretical



Mendel Generic Linux (SL6) 750 Nodes FDR IB Supports HEP, Joint Genome Institute,

Office of Science Materials Genome Project

#### Global Filesystem and HPSS Data Archive







#### **PDSF** Workloads





- serial, high throughput
- Univa Grid Engine batch system
- broad user base (including non-LHC)
- Fair share scheduling projects "buy in" to PDSF and the share tree is adjusted accordingly
- STAR Tier-1
- ALICE Tier-2
- ATLAS Tier-3





### We've Moved





#### ~10 km from Downtown Oakland to LBNL

- Networking 4 x 100 Ge 05.2015
- First Production System part of PDSF – 07.2015
- First unplanned site wide power outage at OSF 08.2015
- Global (GPFS) file systems 09.2015-01.2016
- Cori Phase 1 10.2015
- Second unplanned site wide power outage at OSF 11.2015
- Edison 11.2015
- People 12.2015
- All of PDSF and parent system 02.2016
- HPSS (Tape) 2020?





#### New Building – at LBL









Lessons Learned (?) from the Move



- Takes all one's focus
- Much easier to provide new service/system than migrate an existing service (no surprise here)
- Biggest problems with power
  - Bad circuit breakers
  - Insufficient torque on connections within breaker panels
  - Two site wide scheduled power outages
  - Despite all the tests before the move, trouble with UPS
- Biggest success migrating tens of PB in GPFS from one site to another via live replication without the users noticing (almost)
- Many aspects of the move are serial schedules slip very easily
- Users surprisingly understanding





## What remains at OSF?



- 1. ~60 old computes
- 2. specialized servers
  - xrootd redirectors
  - VO boxes
  - MySQL databases (STAR)
- 3. batch servers (Grid Engine)
- 4. Storage for specific experiments (~426 TB)

- 1. retired at any time
- 2. VMs at CRT ready for migration
- 3. new server and storage hardware for batch system at CRT
- 4. data migration
  - to NGF
  - physical move (half rack)



#### What else is new?



- Use of RDMA rather than IPoIB to access global GPFS file systems

   factor of 3 bandwidth improvement
- ~ 1 PB of storage deployed for ALICE xROOTd storage first EOS at NERSC
- Retirement of 18 PDSF specific GPFS file systems leaving 3
- Our PDSF consultant (app management & support) left in July 2015 and no replacement – position cut in 02.2016 – then brought back – and should have one in 05.2016
- Greater usage of NERSC "global" GPFS file system
  - Changing DTN endpoint for ATLAS grid file transfers
  - Centralized management, greater ability to grow
- Greater use of container technologies on our Cray systems shifter – providing an SL 6 chroot env and a dynamic mirror of necessary cvmfs file systems
- With retirement of legacy systems at old site, only 3 similar hardware flavors remaining







- ½ rack 3 NetApp E5560 60 \* 6 TB drives per tray
   2 JBODs, 1 with dual redundant raid controllers
- 3 FSTs, 1 MGM/MQ server
- Given perl install script we don't have yum/rpm on the servers (due to our provisioning model) – with some changes – was able to install EOS atop SL 6.7 image
- Documentation somewhat scattered and has some hidden assumptions (ELI5 would be good)
- ALICE developers very responsive and helpful





## + EOS architecture

EOS organizes filesystems in views by spaces, nodes, groups and filesystems. By design there can be an arbitrary number of spaces. There should be at least as many groups as filesystems per node.









- cfengine -> ansible for configuration management
- investigate using slurm for batch system
- Start testing SL 7 compute images
- Start migration from chos (simple kernel module providing chroot environment) to shifter (more flexible container environment)
- Lots of thought about increased data from ATLAS experiment in the coming years from Pb-Pb





#### A New Focus – We Host Fall 2016 HEPiX











#### National Energy Research Scientific Computing Center



