



# BNL RACF Site Report

HEPIX SPRING 2015 – OXFORD, UK

WILLIAM STRECKER-KELLOGG [WILLSK@BNL.GOV](mailto:WILLSK@BNL.GOV)

# RHIC/ATLAS Computing Facility

2

- ▶ Brookhaven National Lab
  - ▶ About 70 miles from NYC
  - ▶ Hosts RHIC
- ▶ Provides computing services for both RHIC experiments and serves as ATLAS Tier-1 for the US.
  - ▶ Around 46kCpus of computing
- ▶ Supports smaller experiments—LSST Dayabay, LBNE, EIC, etc...
  - ▶ <1kCpu of computing



# Linux Farm Updates

- ▶ RHIC purchased 160 Ivybridge servers from Dell
  - ▶ 62 kilo-HEPSpec06
  - ▶ 6.7PB Disk Storage
    - ▶ 12x4 TB Drives in RAID5
- ▶ USATLAS Bought 77 Ivybridge nodes
  - ▶ 30 KHS06
  - ▶ Disk-light, used for scratch space only
- ▶ Evaluating Haswell—See Tony Wong's presentation on remote evaluation



# Infrastructure Updates

- ▶ New CRAC unit providing 30 tons of additional cooling to BCF
- ▶ Future Plans
  - ▶ Replace aging PDU and upgrade APC batteries
  - ▶ Provide additional 200kW of UPS power to CDCE
  - ▶ Refurbish 40-year-old air handlers in basement of BCF
- ▶ Formulating proposal for upgrading building 725 (old NSLS) into a datacenter
  - ▶ Little space left in our rooms, additional space needed to support NSLS-II and other research activities around the lab

# Evaluating Software

5

- ▶ Looking at Docker and batch-system support for containers
  - ▶ Useful to support non-native (linux-based) workloads on our dedicated systems
- ▶ Evaluation of SL7
  - ▶ Support in our config-management system is ready
  - ▶ Already used for new infrastructure (DNS hosts / NTP servers, etc...)
  - ▶ No plans to upgrade Farm soon



# HPSS Tape Storage

6

- ▶ Over 60PB in 53K Tapes
- ▶ Writing new RHIC data to LTO-6 tapes, 16 drives
  - ▶ ATLAS moving to LTO-6 soon
- ▶ Migration from lower-capacity cartridges ongoing (LTO-3→LTO-5)
- ▶ Completed testing for Dual-copy with mixed T10KD and LTO-6 cartridges
  - ▶ Copying old tapes halted during RHIC run



# HPSS Tape Storage cont...

- ▶ Performance Measures
  - ▶ New LTO-6 drives gets 133MB/s writes
  - ▶ T10KD drives writes 244MB/s
- ▶ New bug discovered while packing tapes
  - ▶ Intensive repacking leads to issues with monitoring that has disrupted the production system
- ▶ Upgrade planned by end of 2015, for bugfixes and RHEL7 client support

# Central Storage—NFS

- ▶ 2 Hitachi HNAS-4100 heads are in production
  - ▶ Have purchased an additional 2 HNAS-4100 heads
- ▶ Phased out 4 Titan-3200, with an additional 2 Titans to remain in production through this year
- ▶ Mercury cluster for ATLAS is being phased out
  - ▶ Goal of *no dependency on local NFS storage for ATLAS jobs*



# Central Storage—GPFS

- ▶ Fully puppet-deployable on the client
  - ▶ Still needs server configuration—could be solved with exported resources
- ▶ Updated client and server to 3.5.0-18, no problems
- ▶ Sustained spikes observed up to 24 GB/s without problems for the clients
- ▶ Had to put Cgroup throttle on xrootd's /dev/sda usage because GPFS commands had sufficient latency to drop the node from the cluster

# Distributed Storage Updates

10

- ▶ PHENIX dCache and STAR XrootD now over 7.5PB each
- ▶ “Raw” data on PHENIX dCache expanded to 3PB
  - ▶ Used for tape-free “fast” production runs
  - ▶ A portion is used as a repository for special PHENIX jobs running on OSG resources
- ▶ Scratch space for PHENIX production moved from dCache into GPFS
  - ▶ Performing very well

# ATLAS Storage

11

- ▶ dCache upgraded to version 2-10-20
- ▶ ATLAS xrootd n2n plugin on xrootd door
  - ▶ Providing ATLAS N2N and redirection to remote redirector when a file is not found
- ▶ Xrootd Versions
  - ▶ 4.1.1 reverse proxy
  - ▶ 4.x.y (latest git for bugfixes) forward proxy
- ▶ FTS 3.2.32
- ▶ ATLAS jobs
  - ▶ Moving from copy-to-scratch to direct-IO via xrootd door

# CEPH Developments

12



- ▶ Storage backend extended to 85 disk arrays
  - ▶ 3.7k 1Tb HDDs, 3x replication for a usable capacity of 1Pb
- ▶ Head node network can now handle 50GB/s
- ▶ Running 2 separate clusters—main ATLAS prod (60%) and federated CEPH(40%)
  - ▶ New cluster components based on kernel 3.19.1 and CEPH 0.87.1
  - ▶ Tests of CephFS and mounted RBD volumes were performed with 70x1GbE attached clients
- ▶ Long-term CephFS stability tests with kernel 3.15.5 and Ceph v0.80.1 show stability across a 6 month period.
- ▶ Please refer to the dedicated talk by Ofer Rind

# S3 CEPH Throughput Test

13

## Destination CEPH BNL S3 bucket accessed through RADOS GWs(2)

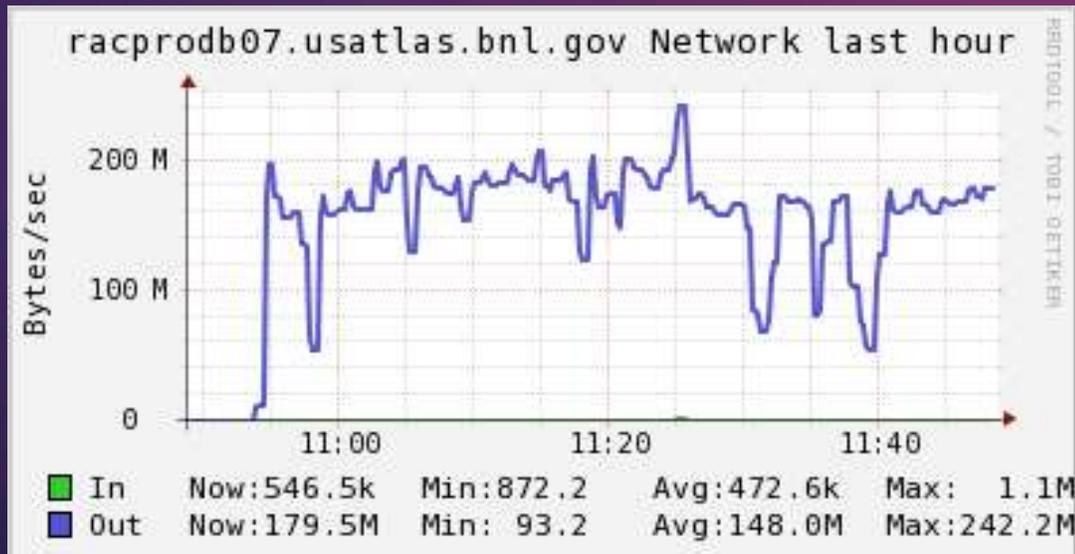
- ▶ Two copy jobs sequentially started 10 minutes apart
- ▶ Each job copied 202 x 3.6GB files
- ▶ 128 transfers/job in parallel allowed
- ▶ 242MB/s maximum network throughput observed



Proxy was populated with 202 x 3.6GB files to run a custom job to access CEPH via python/Boto

## Test server specifications

- ▶ 47GB memory total
- ▶ 16 processors
- ▶ 2x10Gb/s LACP Network connectivity
- ▶ Data stored in s DS3500 backend storage
- ▶ 17 SAS disks /15kRPM configured in a
- ▶ RAID 0 layout



# Configuration Management

14

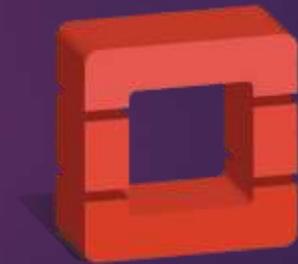


- ▶ Puppet 3.7.4
  - ▶ Migrated away from passenger/apache to new puppetserver
  - ▶ Faster catalog compilation time, fewer errors seen
  - ▶ Work on Jenkins automated integration testing
  - ▶ See William Strecker-Kellogg's talk on Puppet later in the week

# OpenStack in Production

15

- ▶ Icehouse cluster
- ▶ 47 hosts (16CPU, 32Gb RAM, 750Gb Disk)
  - ▶ Second equivalent test cluster slated for Juno
- ▶ Used internally for ATLAS, 3 external tenants
  - ▶ ATLAS Tier-3
  - ▶ BNL Biology Group
  - ▶ BNL Computer Sciences Group



openstack™  
CLOUD SOFTWARE

# Amazon Pilot Project

16

- ▶ Since Sept. 2014, major project to run ATLAS at full scale on EC2
  - ▶ Compute: 3 regions, ~5k node test in November
  - ▶ Storage: EC2-based SE (SRM + s3fs)
  - ▶ ATLAS application customizations
    - ▶ S3-enabled stage in/out
    - ▶ Event service (checkpointing) to leverage spot pricing
- ▶ See talks later in the week by John Hover (BNL) and Dario Rivera (Amazon)

# The End

Questions? Comments?  
Thank You!

