

Integration & Optimization of Computing within BNL workload management system

Iris Wu
ISGC 2018, Taipei

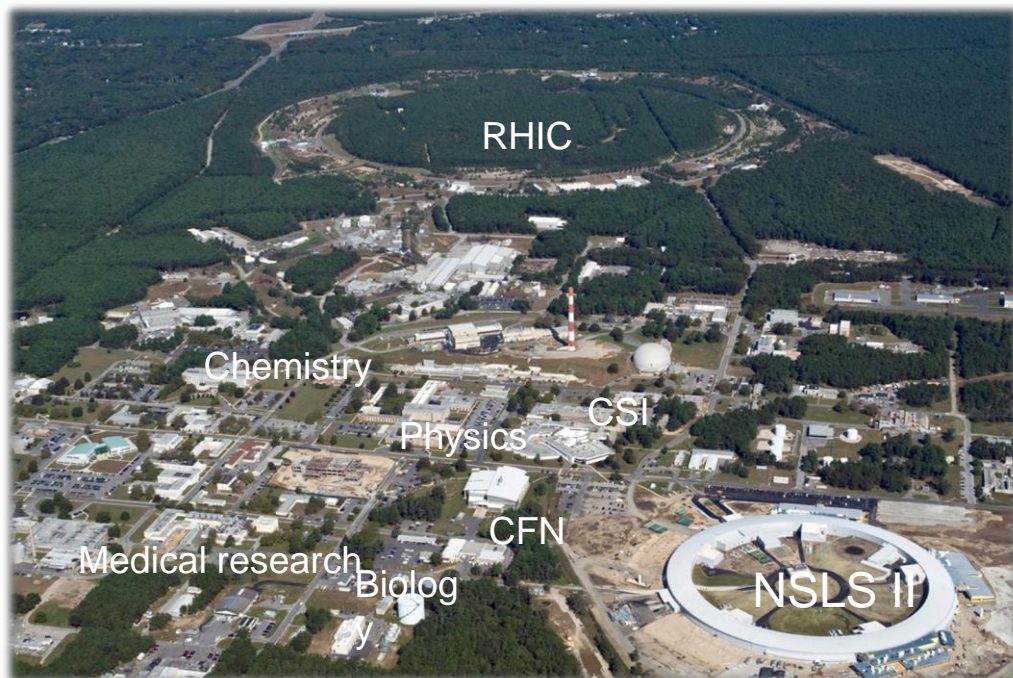
70 YEARS OF
DISCOVERY

A CENTURY OF SERVICE



Scientific Data & Computing Center(SDCC)

- Located at Brookhaven National Laboratory on Long Island, New York
- Provides full computing service for:
 - RHIC experiments – STAR, PHEINX
 - Tier 0 center
 - RHIC RUN 2018 started
 - ATLAS
 - Largest ATLAS Tier 1 center worldwide
 - Delivers ~ 200 TB/day
 - Small groups: LSST, Daya Bay etc.
 - Belle II
 - Tier 1 center
 - Migration of service from PNNL to BNL by Sept 30, 2018
 - Science Domains within BNL other than HEP/NP

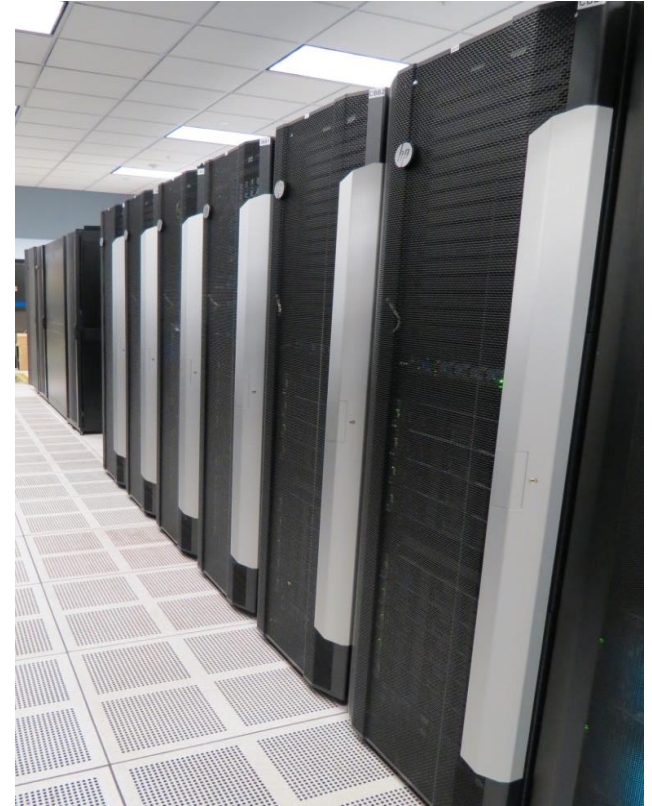


Integration & Optimization of Computing

- HPC/HTC
- Storage
- Network
- Configuration Management

HPC Clusters

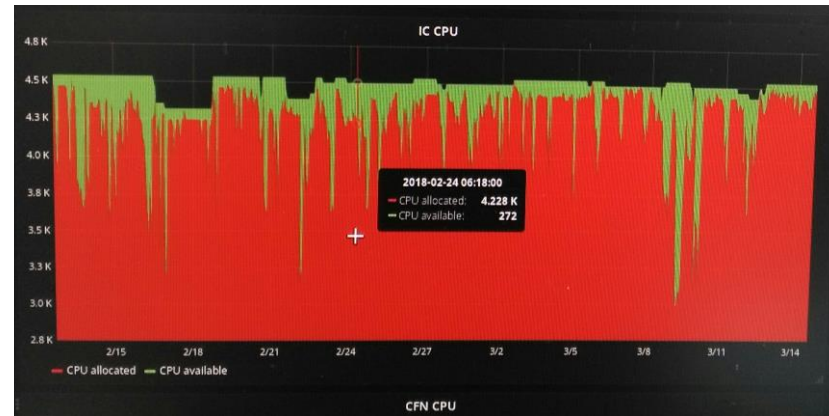
- SDCC has a 30k HS06 (and growing) cluster available to HTC/HPC
- ATLAS and PHENIX have allocations on HPC resources
- Acquired Institutional Cluster (IC) and KNL-based cluster for HPC-based projects
 - Current stakeholders
 - IC : CFN, MS, LQCD
 - KNL : LQCD
- Shared storage
 - 1 PB of GPFS storage with up to 42 GB/s I/O bandwidth capability
 - Seagate G200 appliance connected via infiniband



IC & KNL Clusters

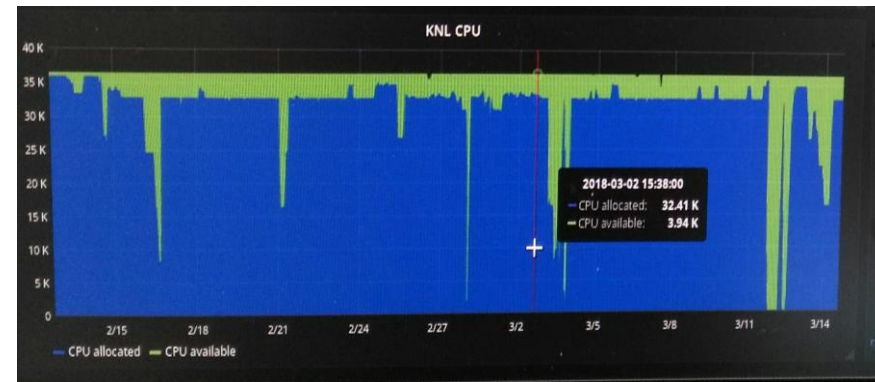
- Institutional Cluster(IC)

- 108 compute nodes
 - Dual Xeon Broadwell (E5-2695 v4) CPU
- 2 head nodes
 - Two NVidia K80 GPU's
- SAS based local storage / 256 GB RAM
- Dual-rail Infiniband EDR interconnect
- Added 18 compute nodes (to be expanded to 108)
 - Two NVidia P100 GPU's



- KNL-based Cluster(KNL)

- 142 compute nodes + 2 submit nodes
- Single Intel Xeon Phi CPU 7230
- SSD-based local storage /192GB RAM
- Dual-rail Omni-Path (OPA) interconnect



Skylake Cluster

- To be delivered to BNL April 2018
- Single-rail Infiniband EDR interconnect
- Configuration of computer nodes
 - 64 dual-socket Xeon skylake Cluster 6150(Gold) with 72 logical cores
 - 4 x 4 TB SATA drives
 - 192 GB of RAM



Linux Farm

- New machines in production since January/February
 - RHIC – 34 Dell PowerEdge R740xd servers, 2448 job Slots, ~35 kHS06
 - ATLAS – 90 Dell PowerEdge R640 Servers, 6480 Job Slots, ~93 kHS06
 - Total 68K logical cores now

- Belle II
 - 10 new machines online
 - same specs of R640 order of ATLAS
 - 56 repurposed Dell PowerEdge R620 servers in production capacity ~28 kHS06
 - 2 Xeon E5-2660 (Sandy Bridge) 2.2 GHz CPUs (32 logical cores total)
 - 1 Gbps Ethernet
 - Systems built with SL7, jobs running in SL6 Singularity containers
 - New machines delivered soon (April 2018)
 - 27 dual-socket Xeon Skylake 6150 (Gold) with 72 logical cores, 28k HS06

HTC/HPC integration



- Developed at the University of Wisconsin
- Highly scalable in handling large number of disparate jobs
- Support FIFO or best fit scheduling only
- Limitations
 - Environment where different users/groups pay for shares of cluster and opportunistic usage is permitted

- Primarily developed by SchedMD
- Fair share scheduling policies for parallel/multi-node jobs
- Partitioning of resources and advanced reservation
- Not capable of handling large number of running/queued jobs

HTC/HPC integration

- HTC ----→ HPC

- A mechanism to submit HTCondor jobs to SLURM Batch system
- Allow users to opportunistically utilize HPC/Slurm resources
- Transparent to users

- HPC ----→ HTC

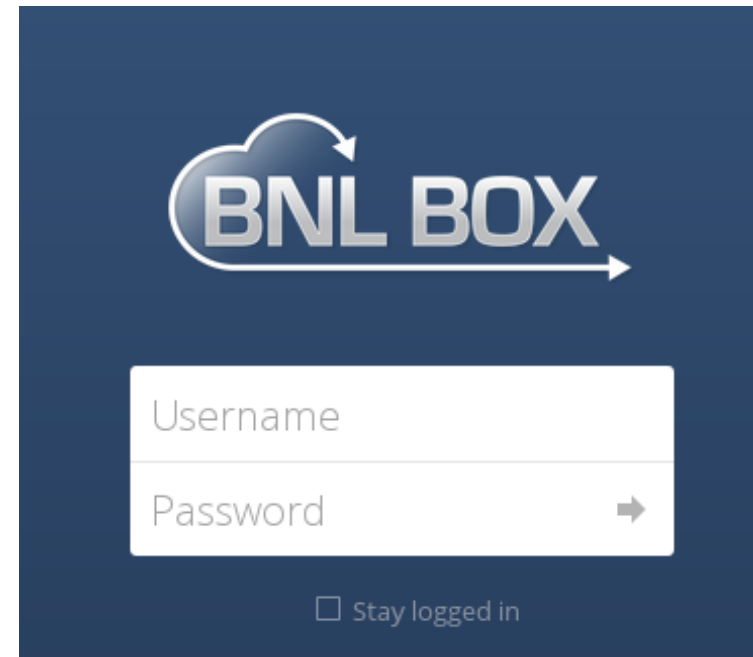
- HPC applications generally require low-latency node to run efficiently
- Possibility
- Users need to manually submit such jobs to HTCondor

Distributed Storage

- ATLAS dCache
 - Version 3.0
 - 17.5 PB disk storage
 - JBOD, Hitachi storage
 - Enforces quota for US Atlas Tier 3 users
 - Resilience dCache
 - Supported Protocols
 - SRM
 - HTTPS/Webdav
 - NFS
 - Direct I/O access (NFSv4.1)
 - xRootD transfer
 - HSM feature
 - Non-Blocking HSM
 - Retry pftp within same request
- Belle II dCache
 - version 3.0
 - 1.7 PB disk storage
- Simon's dCache
 - version 2.0
 - Ceph storage
- PHENIX dCache
 - version 1.9 (in the progress of upgrade)

BNL Box

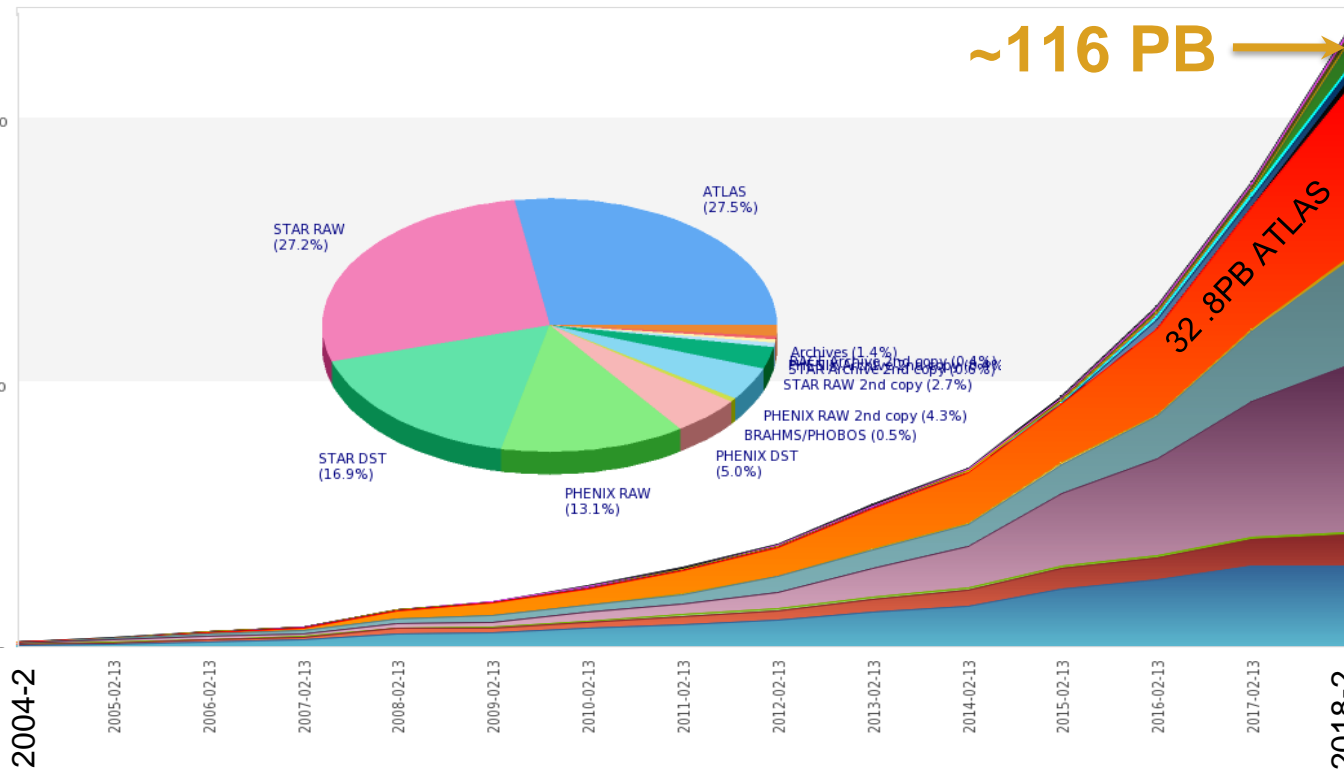
- Store and access data on various platforms
- Using Owncloud software
 - Quota for each users
 - Users can share data
 - Configuration of Synchronization
- CephFS storage
 - 5.1 PB, Luminous 12.2 release
 - 40 Gbps bandwidth
- Stream access
 - WebDAV
 - xRootd



Mass storage- HPSS

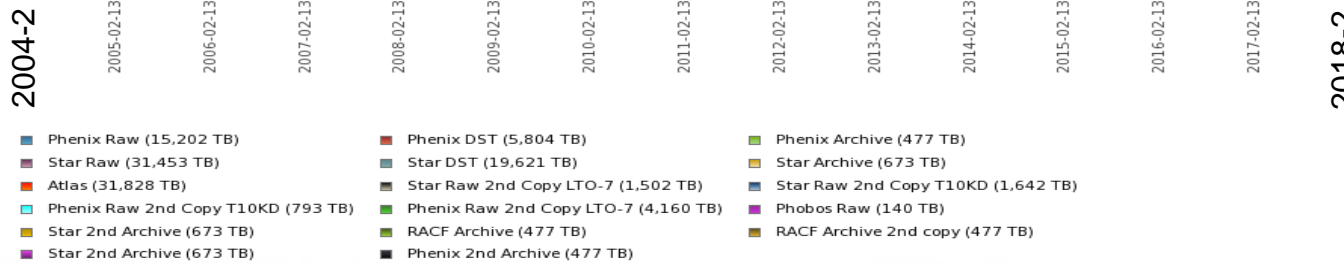
HPSS Data Yearly Growth

Date: [2004-02-13 - 2018-02-13], Total: 113.4 PB
 All Archive Storage Classes are counted as single copy



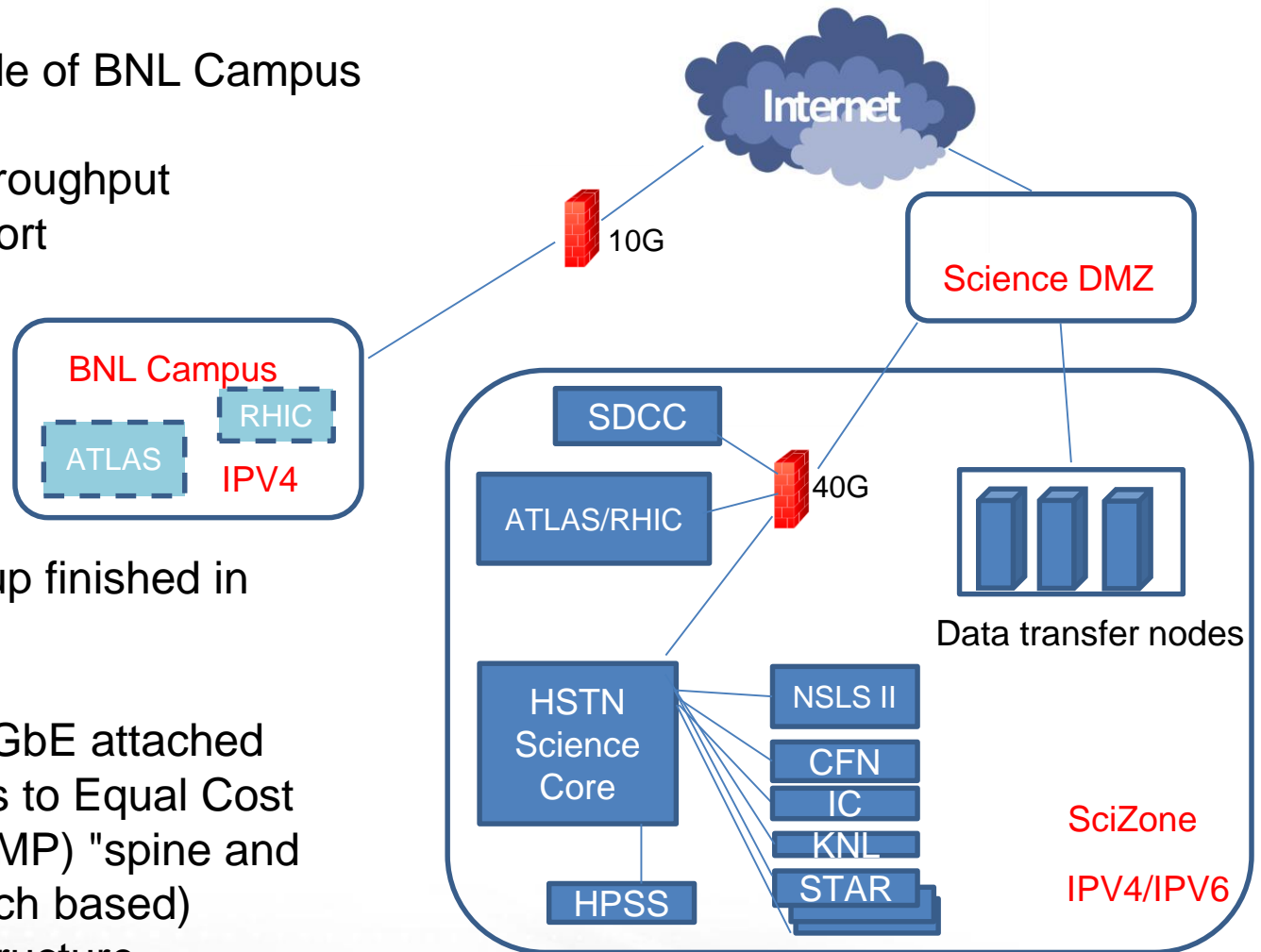
- **Write to tape:**
- Big files
- File family

- **Read from tape:**
- Round robin gateways
- Minimize tape mount



Network Re-configuration

- Migrated outside of BNL Campus network
 - Firewall throughput
 - IPV6 support



- SciCore build up finished in February
- Migration of 1 GbE attached compute nodes to Equal Cost Multi Path (ECMP) "spine and leaf" (ToR switch based) network infrastructure

Configuration Management

- Puppet
 - Version 3.7, evaluation of new version
 - Heavily usage of puppet in SDCC
 - Using git-subtree to share a common set of modules
 - Can import/export a common subdirectory containing shared Puppet modules into a user's gitPuppet repo
 - Jenkins automated integration
 - Migrate Cobbler templates to Foreman



Thank You!

Information here compiled with help from the entire SDCC staff.