



Scientific Computing and Data Facility

(Formerly the SDCC, originally RCF)

Shigeki Misawa
May 8, 2025

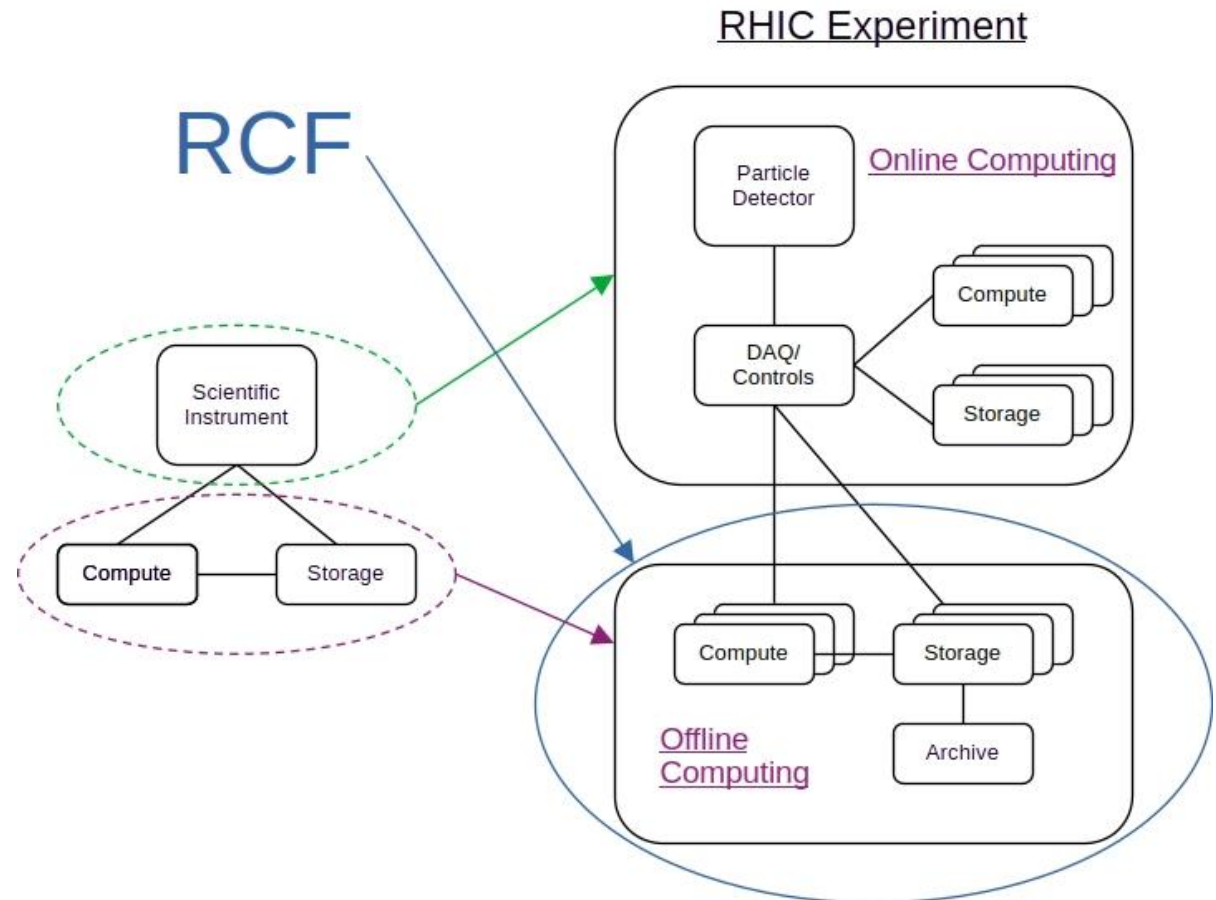


RHIC Computing Facility (RCF)

- RCF Mission: Develop and operate computing and storage infrastructure for RHIC experiments
 - Archive data streaming from the experiment(s)
 - Build/operate dedicated, large scale, high throughput (HTC) compute farm
 - Predominantly large scale, embarrassingly parallel, “batch” workloads
 - Limited, small scale “user analysis”
 - Build/operate large capacity online/nearline storage systems to support computing
 - Dedicated, single source funding to develop the above capabilities.

RHIC Experiments (Now)

- More data, faster
 - 100's of PB
 - 10's GB/sec
- More storage
 - ~100 PB disk
 - ~ 300PB tape
- More compute
 - 1.4K nodes/150K cores
- More people
 - 100's of collaborators per experiment

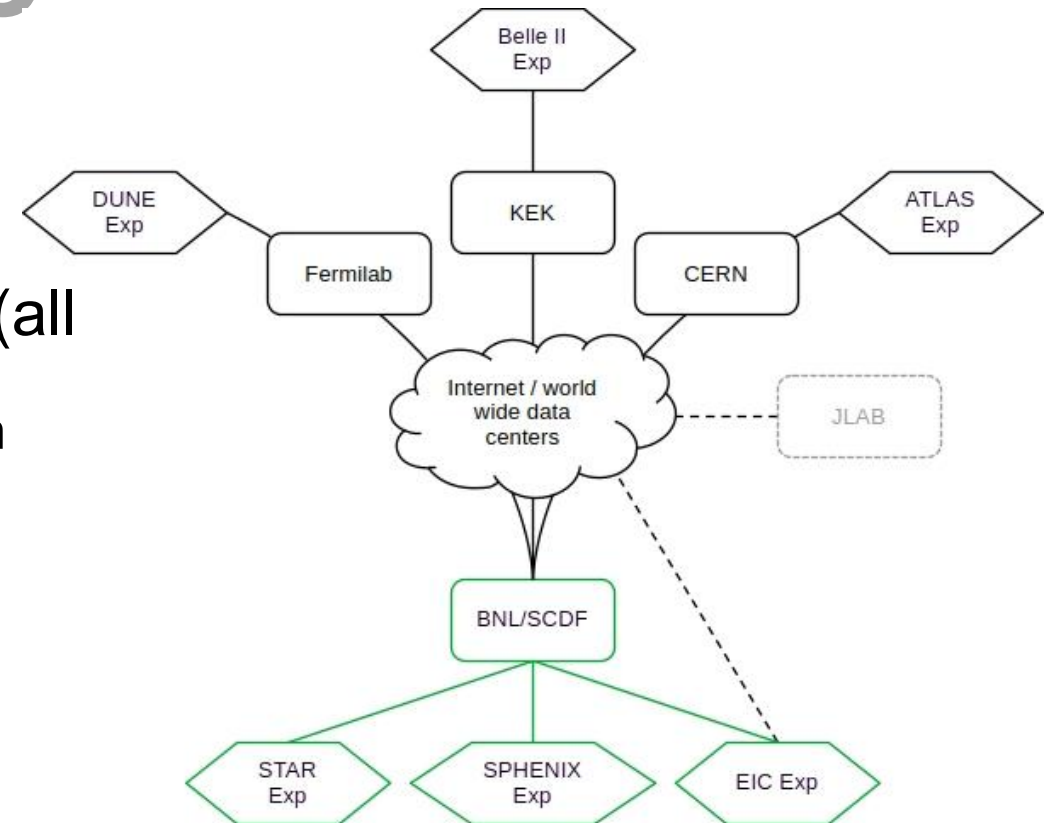


RCF → An Unconventional Facility

- NOT a “full service” computing center
 - Limited set of services – Essential for sustainability (critical mass)
 - Homogeneous, large scale resources – Achieves “economies of scale”
 - Tailored to support “big data”, “big science” experiment communities
 - Effectively a “bespoke” computing facility
 - No real concept of shared or common, re-allocatable resources
 - Not like a leadership class HPC center
- Focused on Operations
 - Virtually all user support delegated to support staff within experiment communities

RCF → RACF → SDCC

- Added support for “adjacent” communities
 - ATLAS, sPHENIX, Belle II, DUNE, EIC (all NP/HEP)
 - Significant overlap in communities resulting in “cross pollination”
 - Virtually identical computing models
 - Similar or identical tools/services
 - Similar support requirements
- **Funding continues to be program specific**
 - Multi-decade experiment lifetimes



RCF → RACF → SDCC

- New experiments located around the world, not limited to BNL
- Added interactions with other major scientific data centers, located worldwide, supporting the new communities.
 - Experiment specific “computing grids”
 - Data center interoperability required for major services
- More primary services
 - Global data and workflow management in collaboration with NPPS (Nuclear and Particle Physics Software Group)
 - Federated identity management

RCF → RACF → SDCC (cont'd)

- Some dilution of critical mass/economies of scale, e.g.
 - Evolutionary changes in software technologies
 - Additional effort deploying new technologies embedded in existing services (e.g. transition from PKI to OIDC token based authorization)
 - Integration of new community software and work processes with facility services
 - Increased diversity in “secondary” services
 - e.g. Ansible vs Puppet vs Chef, Mattermost vs Discord
 - Additional “organizational” overhead
 - Work coordination, accounting, communication overhead
- Still predominantly a “big data”, “big science”, high throughput computing facility, but now with a global reach

RACF → SDCC → SCDF

- Support for new, “non-adjacent” communities
 - HPC and smaller “big data” groups
 - Different computing models
 - Limited commonality in tools and services with traditional SDCC clients
 - Different support requirements
 - Minimal to no overlap in communities (existing or new)
- HPC support through work with USQCD and CFN Theory and Computation group
- Small “big data” group support via NSLS-II and CFN Electron Microscopy
- Small group support with the Cosmology and Astrophysics Group

RACF → SDCC → SCDF (cont'd)

- Funding continues to be program specific
- Viable funding model for smaller groups problematic
 - Smaller and more “transient” clients
 - Difficult to achieve critical mass/economies of scale to support services
 - Participation of multiple groups required for sustainability/affordability
 - Financing upfront capital costs problematic
 - Stable, long term funding also an issue
 - Strong impact on staffing
- Funding possibilities
 - BNL LDRD or Program Development (PD) funding
 - Addition of funding for computing in future grant proposals

Key Facility Takeaways

- At this time the SCDF is a “lean”, operations focused facility
 - User/software support outsourced to customers
- “Programmatic” funding model
 - Configured for large, long lifetime collaborations with customers
 - Resources and services are dedicated to the funding program
 - This include FTE effort
- Critical mass/economies of scale in a selected set of capabilities
- R&D problematic for the above three reasons
 - Extension of existing services, as is, to new customers most efficient
 - Deployment of new or highly customized services difficult.
- However, the SCDF is open to discussions on a sustainable model for expanding SCDF’s support of the research community

Comments : R&D

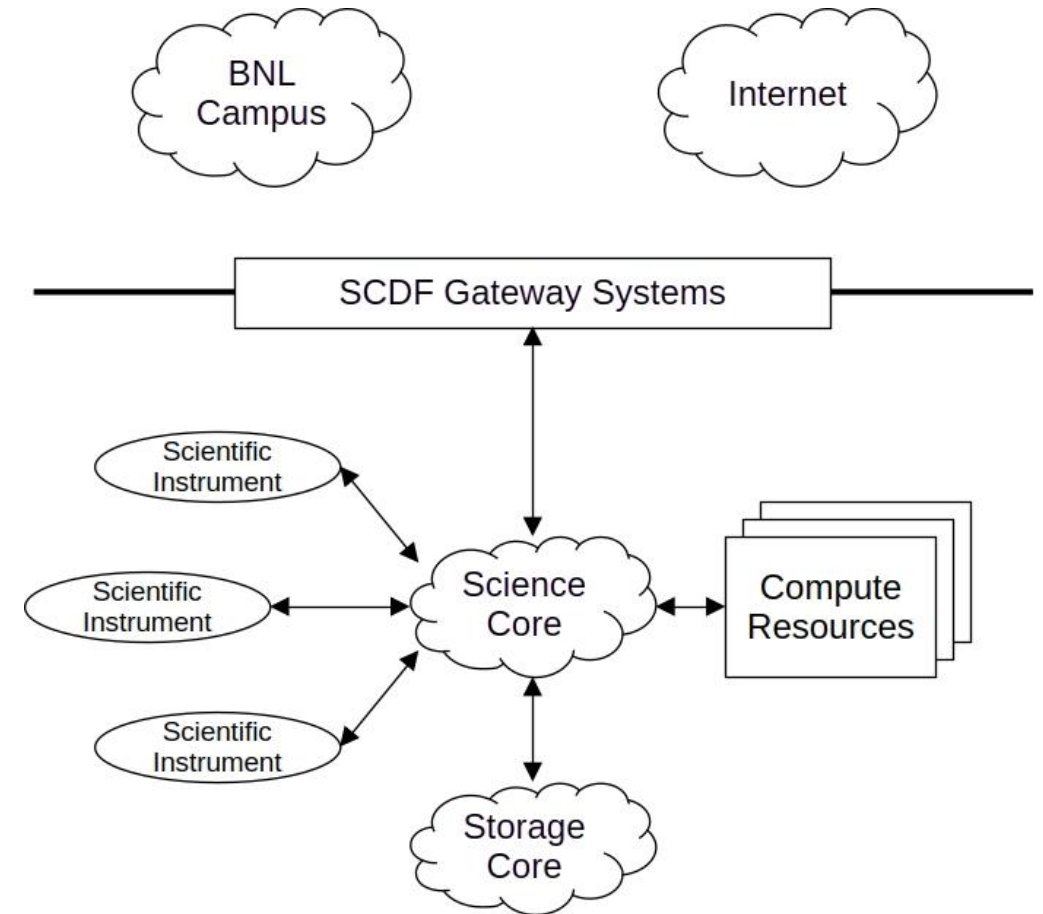
- R&D and “one off” systems are “high touch” services
 - Support costs are likely to be high with periods of substantial effort being episodic.
- Funding for acquisition of systems and support staff are program specific
 - No “SCDF” funding to finance purchases of test systems
 - No process in place to request access to “exotic” systems purchased by other programs
- Limited colocation service (equipment hosting) is available
 - But not in the SCDF production data centers
 - In BNL campus network, not in SCDF network.
- SCDF has vested interest in investigating new technologies that may be of use to the broader BNL community
 - Framework needs to be developed to encourage collaboration among all SCDF stakeholders on R&D

SCDF by the numbers - Data Centers

- Bldg 725 data center (Rm 1-301)
 - N+1 redundancy
 - Power and cooling self sufficient
 - not dependent on utility power or BNL central chiller plant
 - PUE 1.2 to 1.4
- Bldg 515 data center (CDCE/Rm 1-4K)
 - Power – self sufficient
 - **Dependent on central chiller and utility power for cooling**
 - PUE Significantly higher than 1.4
- Full cost of each data center must be borne by the data center occupants, not just the occupied areas.
 - Worst case, single occupant with one rack in CDCE pays all space charges for the room
 - **Program space charge will change as other programs add or remove racks from the data center**
 - Bldg 725 occupants pay for cost of supporting electrical and mechanical rooms.
 - Effectively triples cost of space charges for racks in Bldg 725.

SCDF : A Standalone Facility

- SCDF is not “in” the BNL campus network
 - From BNL, SCDF is part of the general Internet
 - Similarly, from the SCDF, BNL is part of the general Internet
- SCDF network
 - High Throughput Science Network (HTSN)

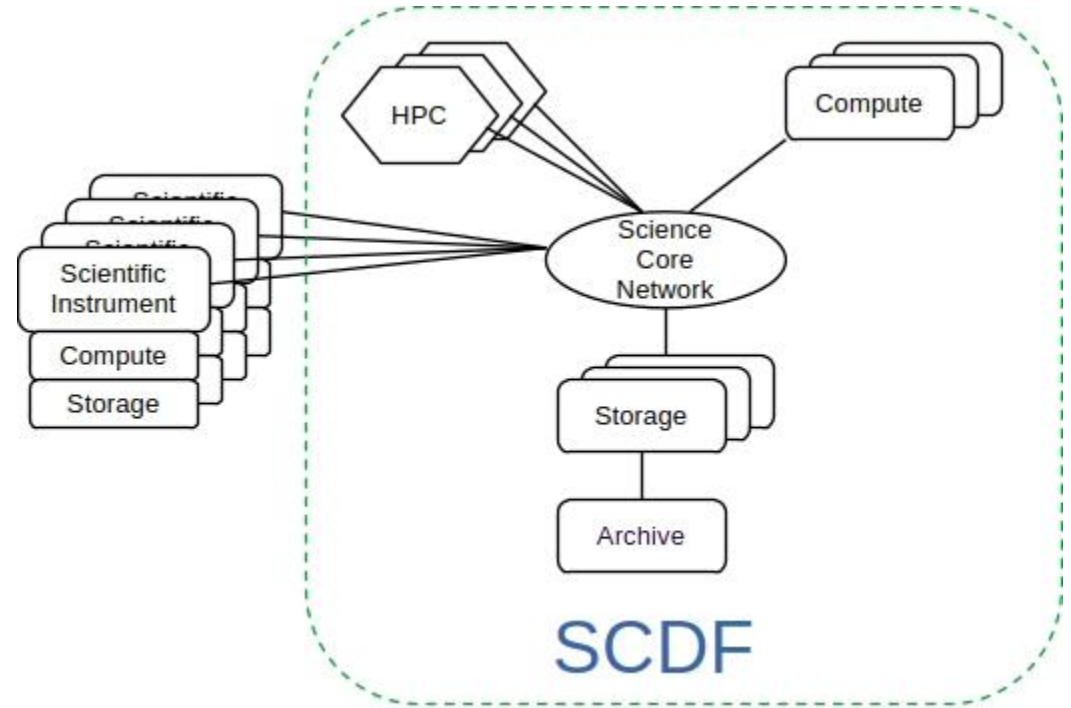


High Throughput Science Network (HTSN)

- Dedicated scientific data network for high bandwidth connectivity
 - Science DMZ
 - Landing point for high bandwidth data transfers to the WAN.
 - 1.6Tbps aggregate bandwidth to the WAN
 - Science Core - HTSN Tbps routing fabric
 - Storage Core – Network hosting point for data storage
- Infrastructure can be scaled to meet project needs

Science Core Connectivity

- SCDF internal connectivity
 - HPC clusters
 - HTC compute farm
 - Disk storage systems
 - Tape storage systems
- SCDF “external” connectivity
 - sPHENIX experiment
 - STAR experiment
 - CFN electron microscope
 - NSLS-II



SCDF: Equipment Life Cycle Policies

- Equipment life cycles (by default)
 - SCDF network infrastructure – Up to 7 year (End of Support date)
 - Top of rack switches for compute nodes – Until end of support contract or when nodes in rack retire, whichever is shorter.
 - All network equipment covered by vendor support contracts
 - Compute equipment – Typically length of support contract
 - Storage – 5 year life cycle w/ support
 - Tape Equipment - covered by vendor support contracts
 - Tape libraries - Depends on availability of vendor support, but may be shortened by cost
 - Tape drives - Primarily driven by media migration schedule and cost of vendor support
 - Tape media – Migration to new media typically every 2 generations.

SCDF Core Capabilities

- Computing
 - HTC – Ethernet connected compute nodes
 - HPC – Infiniband connected GPU/CPU nodes
- High Throughput Science Network (HTSN)
 - LAN (to experimental apparatus) and WAN
- Storage
 - Parallel file systems
 - Bulk online storage
 - Archive/nearline tape storage
- Data Management and Distribution (local and wide-area) *
 - But outside of Globus, limited to NP/HEP domain specific systems.
- OpenShift virtualization/containerization

□ With NPPS (Nuclear and Particle Physics Software Group)

Questions ?

Backup Slides:

Sample of Managed Resources

- 2.6K node/230K core high throughput compute farm for NPP
- Four HPC systems – one for each of the following; LQCD, CSI, CFN and NSLS-II
- Spectrum of disk storage systems (for multiple programs)
 - Scale out Lustre systems (~94PB)
 - dCache disk storage systems (~200PB)
 - Flash based NFS systems
 - GPFS disk storage systems
- Archive/nearline tape storage (~ 350PB) primarily used by NPP and CSI
- Terabit network infrastructure (NPP/NSLS-II/CFN)
 - 25/100/400 GbE capable

SCDF: Additional Services

- BNLBox - File hosting service (like Dropbox)
- Jupyter Notebook - Web based data analysis platform
- Website development
- Document Repositories - Development of Invenio based digital repositories
- Mattermost chat service - Opensource alternative to Slack.
- CVMFS - Software distribution platform
- Rucio - Global data management and distribution service
- PanDA - Distributed workload management service