# dCache at BNL

**Carlos Fernando Gamboa on behalf of the BNL Storage / dCache Team**
**Brookhaven National Laboratory**

**17th International dCache workshop  2023, Berlin, Germany 2023**

1

# Outline

- Overview to dCache based storage services
- Toward an improved dCache service
- Challenges and future work

# Storage Services at BNL SDCC

- BNL SDCC supports different storage services for a variety of Scientific Communities (SC) like NSLSII, Nuclear and High Energy Physics
- Diverse storage technologies are used to support the communities: dCache, Lustre and GPFS, please see past HEPIX 2023 BNL site report for specifics
- This talk will concentrate on **dCache storage** technology
  - dCache services for LHC-ATLAS, BELLE2 and DUNE SC store and manage 76PBs (30% DISK) of data distributed in around 122M files
  - Scientific Community data is produced outside BNL:
    - CERN (Switzerland/France),
    - KEK (Japan),
    - Fermilab(IL,US)

Brookhaven
National Laboratory

# Who We Are

dCache Service Application and Operation Support:
Carlos Fernando Gamboa
Vincent Garonne
Qiulan Huang
Matt Snyder (newcomer)


Rob Hancock (Hardware)
Kevin Casela (Hardware)

Brookhaven
National Laboratory

# Experiments External to BNL



**ATLAS**, 66M files and 20PB on disk, 36M files and 50 PB on tape

**BELLE**, 18M files and 2.4 PB on disk, 2M files and 2.7 PB on tape
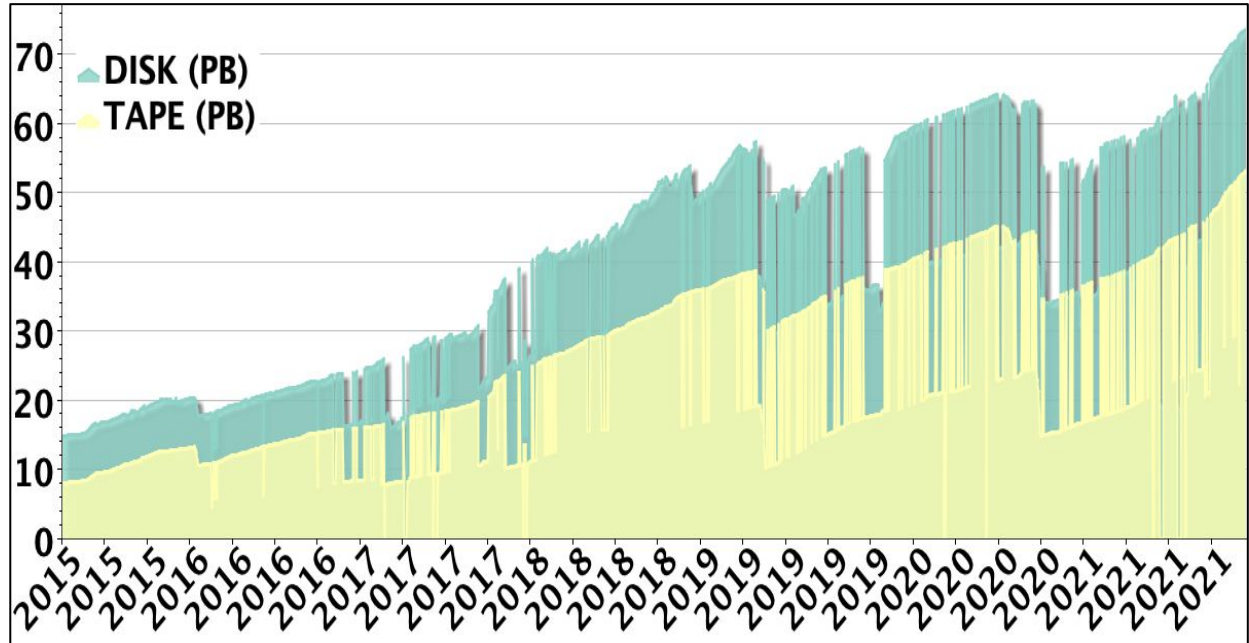
**DUNE**, 500 TB in 200k files on disk

ATLAS SC community driving the storage usage compared to other HEP SC supported at BNL
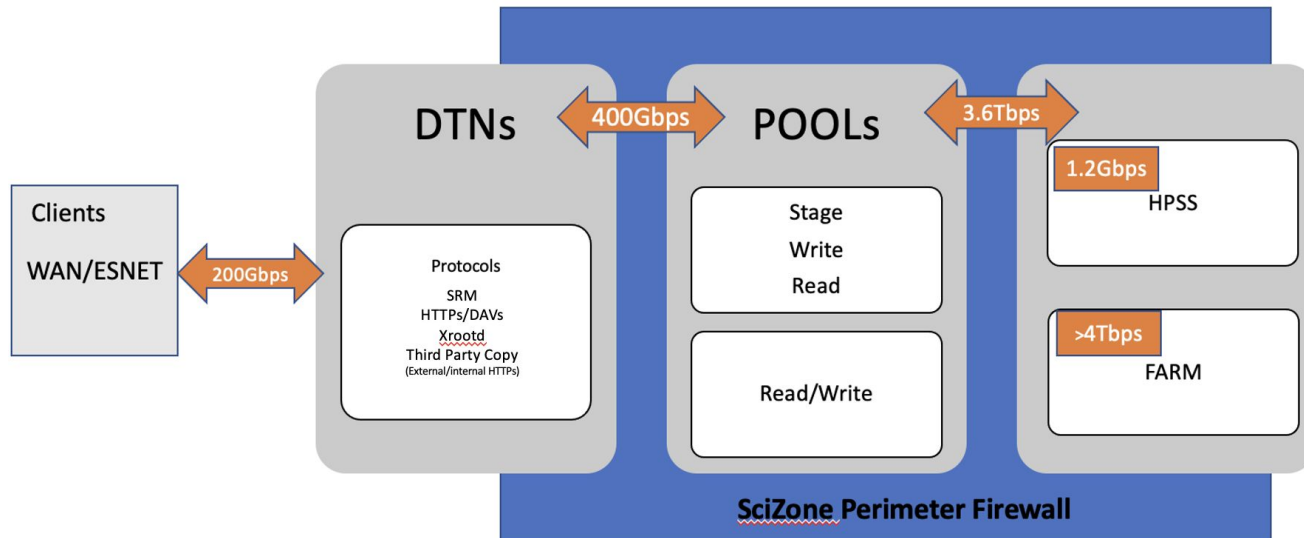
# Evolution of Atlas SC storage

BNL provides more than 70PB of storage and hosts 100M files for ATLAS

We observe a factor of 7 in the past seven years for the total space.



The main challenge coming is HL-LHC and with the simple model of 3 to 4 order of magnitude increase in 10 years from now: 1B files, 700 PB, 300Mhz, 5-7PB/day
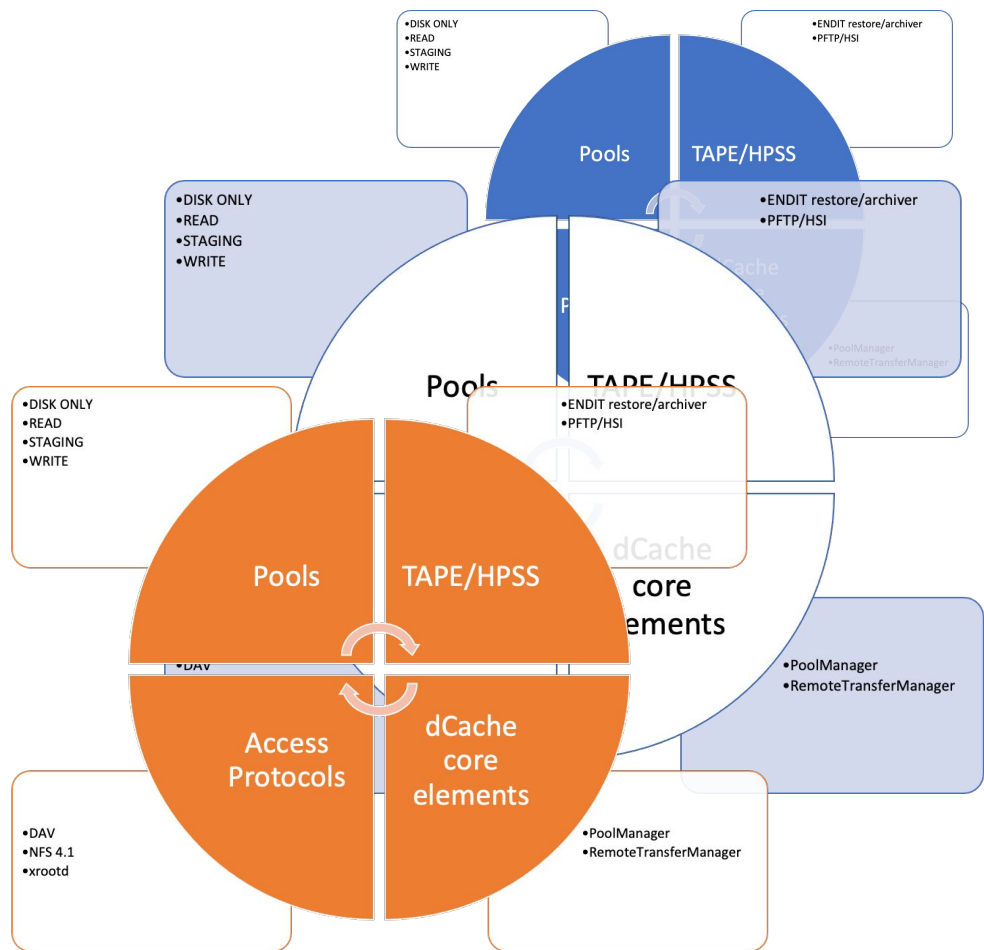
# dCache General Layout (ATLAS)



Comply with BNL cybersecurity policy disaggregation among external and internal resource accessibility

Reference deployment to be used as building block for other SC

Brookhaven
National Laboratory

# dCache instances are isolated per SC

- SC diverge in their requirements
- Procurement and resource control
- Infrastructure supported on physical and virtual Machines

# Towards an Improved dCache Operation

**Areas of work:**

- Enhancing software for interaction among dCache and TAPE HPSS systems
  - ENDIT archiver/retriever
- Improving dCache data access workflows for client access
  - Non firewalled Xrootd client access for write/read
  - DUAL IPv4/IPv6 dCache application stack configuration
- Extending monitoring for dCache operations
- Evolving dCache along with infrastructure

# Improving Software to Interact with dCache and TAPE

**ENDIT  archiver/retriever**

- Previous mechanisms used to instantiate restores from HPSS relied heavily on polling the dCache Poolmanager
- Stability of Poolmanager component at risk when > 100k concurrent requested restores
- Since ENDIT retriever adoption, **no more Poolmanager stability issues were observed**, more than 140k concurrent restore requests without any issue
- Successful adoption of ENDIT retriever permitted the extension of usability for writing interactions to HPSS
  - Allowed consolidate legacy software/code for writing to HPSS

Extended overview covered on this talk

Brookhaven
National Laboratory

# Non Firewalled Xrootd Client Access for Write/Read

Standard xrootd client transfers involve pool redirections among client and dCache service

- Accessibility to clients outside BNL to pools is not permitted

**Support for xrootd in proxy mode released on** dCache 8.2.2

- Proactive functional test work along dCache Developers (Al Rossi et al.)

- First enabled on DUNE dCache to READ/WRITE via xrootd

- Later on successfully integrated on ATLAS dCache instance (8.2.15)

  - Xrootd standalone servers used to front dCache xrootd to provide xrootd external READ (ATLAS) decommissioned

Brookhaven
National Laboratory

# ATLAS DUAL IPv4/IPv6 dCache Stack Configuration

Latest dCache upgrade (8.2.15) permitted to:

- Utilize dual-stack network infrastructure deployed on different components (doors, core, and pools)
- Configure the dCache stack to be able to support client requests on IPv6 and IPv4 in dual networks:
  - dCache doors configured to support different client accessibility
    - Clients internal to BNL LAN supporting only IPv4 or IPv6 (no proxy access)
    - Clients external to BNL proxied access for IPv6 and IPv4
  - **The use IPv6 when transferring data between two dual-stack machines for HTTP-TPC transfers**

Monitoring is key to help identify different data workflows

Brookhaven
National Laboratory

# Monitoring Enhancement

Grafana based monitor using the dCache billing/chimera/srm databases to provide information use in operations

Allows aggregate information from different dCache events by entering the PNFSID (dCache file ID)



Performance of dCache

13

# Monitoring Used in ad-hoc Studies:
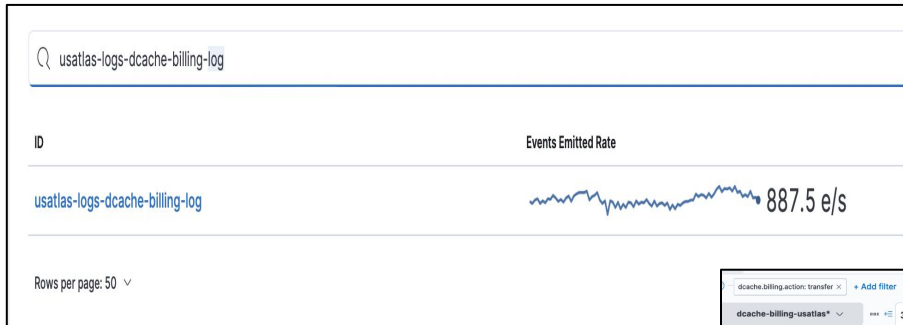
Allowed us to identify areas of improvement for dCache resource access



**Ongoing work to optimize BNL to BNL HTTP-TPC resource data access**

# dCache and ELK Stack
# Started to be Used in Operations

Filebeat / Logstash  pipelines enabled for domain logs and billing logs



ELK use to mine the billing logs with arbitrary queries

# Evolving dCache Along with Infrastructure

- ○ SDCC puppet new infrastructure evolving from puppet 3 to puppet 8
- ○ dCache related puppet modules in principle ported to puppet 8
- ○ New effort in refactoring dCache puppet classes for a multi-instance deployment
- ○ RHEL 7 ~ 2 years for end standard support, new hardware deployment on RHEL 8

| dCache instance | Number of VMs+Physical Hardware(PH) | OS RELEASE | dCache Version | Notes | Pools storage filesystem transition to ZFS from MDRAID |
|---|---|---|---|---|---|
| ATLAS | 84(96%PH) | RHEL 7.8 | 8.2.15 | Hardware for core services to be upgraded end 2023 | 20 servers |
| BELLE2 | 10(100%PH) | **RHEL 8.6 (Core services), Pools (7.8)** | 7.2.19 | Core services hardware recently refreshed (Aug/2022) Upgraded from 6.2.x | |
| DUNE | 5(60%PH) | RHEL 7.8 | 8.2.2 | Resilient manager recently commissioned | 4 Servers |
| Pre-production | 5(100%PH) | RHEL 8 Doors/Other (7.8) | 8.2.18 | WLCG REST API test endpoint Integrated with ATLAS DDM test infrastructure | 1 server |

**Brookhaven**
National Laboratory

# Future Work

- Consolidation of software stack on dCache migration to 8.2.X releases across instances

- Migration of hardware ATLAS dCache to 8.2
    - Hardware refresh cycle→ RHEL 8.2 → Puppet 8
    - Refactorization puppet code for a multi instance dCache deployment
    - Migration of hardware opportunity to consolidate hardware into new datacenter

- Possibilities to enhance monitoring (ELK stack for Billing DB and components log events)

- Participation on HTTP REST Tape API testing

**Brookhaven**
National Laboratory

# In Summary

BNL SDCC is supporting dCache based storage for a diverse of SC

Evolution of the dCache storage features adapted to SC

Priority work will concentrate on:
> Review/optimize TPC data flows for internal transfers (BNL to BNL)
> Improving the orchestration management of dCache software
> Consolidating dCache release levels and OS in different instances

Brookhaven
National Laboratory

# Thank you  /  Danke schön