# BNL ENDIT System: Efficient dCache Interface to HPSS

Zhenping (Jane) Liu

Scientific Data and Computing Center, Brookhaven National Laboratory

@BrookhavenLab

03/21/2024

# The Challenges Leading to BNL ENDIT Development

- **Bottlenecks and Performance Issues:**
  - dCache staging limitations exposed during intense WLCG tape challenges.
  - Out-of-memory errors; non-functional servers during heavy restore requests of 120K or more.
  - Large amounts of duplicated restore requests submitted to HPSS Batch system due to PoolManager restore retry upon pool crash.
- **dCache's default driver interface with HPSS tape storage struggled with scalability.**
  - Scalability significantly limited due to the synchronous nature of this approach and the high resource demands resulting from periodic invocations of an external house-made HSM script for every file being staged.
  - Excessive concurrent PFTP connections from dCache pools caused HPSS PFTP gateway connection problems.

Brookhaven
National Laboratory

# Motivation for ENDIT

- A scalable solution needed to efficiently handle increasing staging demands.
- BNL adapts NDGF's ENDIT system ideas, developing a customized version.
  - Specifically enhances HPSS tape storage interfacing for BNL's requirements.
  - Designed for robust performance in both staging from and migration to HPSS.
- Running in ATLAS dCache production for two years

Brookhaven
National Laboratory

# Components of BNL ENDIT System

**ENDIT Provider:**

- Adapted from NDGF plugin with minor BNL-specific customizations (by Vincent Garonne)

**ENDIT Daemons:**

- **HPSSRetriever Daemon:**
  - Submits stage requests to the HPSS BATCH System.
  - Retrieves files from HPSS via PFTP after HPSS BATCH system processing.
- **HPSSArchiver Daemon:**
  - Responsible for flushing files from dCache tape write pools to HPSS.

**Cron Jobs on HPSSBATCH:**

- One job verifies new staging requests, ensuring only unique and non-duplicated requests are submitted to the BATCH system.
- Another job sends callbacks to the requesting pool(s) upon completion of file processing in the batch system.

**Brookhaven** National Laboratory

# dCache ENDIT Provider plugin

- Adapted from NDGF.

- Minor changes (additional metadata) to accommodate BNL's

- Mechanisms
  - Use predefined file actions within specified directories based on the status of requests.
    - File actions: creating, modifying, and deleting specific files in specific directories
  - WatchService: Watching file events triggered by those file actions
  - Java's and Google Guava's concurrent frameworks: handling file events asynchronously.
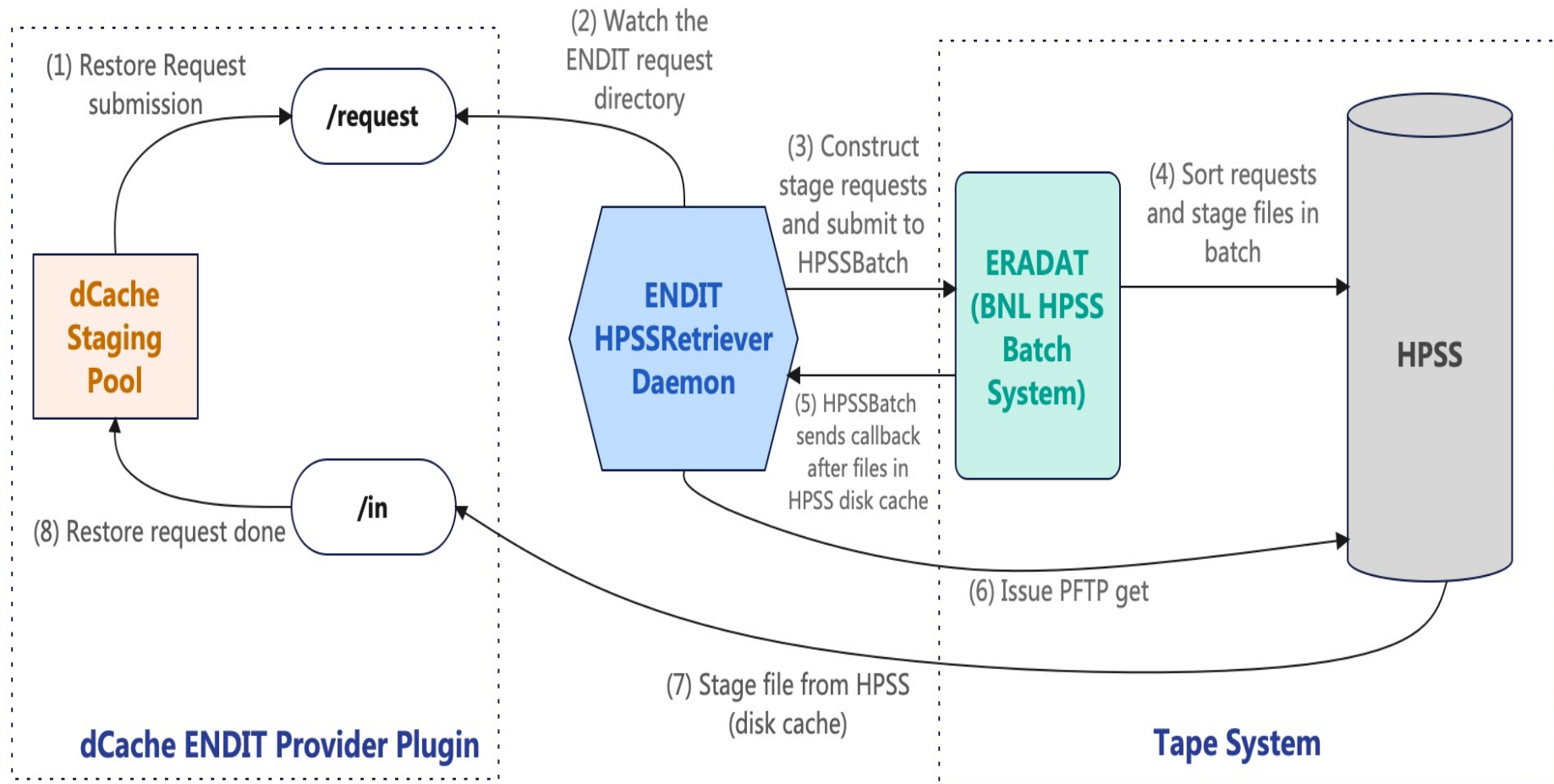
Brookhaven
National Laboratory

# dCache ENDIT Provider plugin (Cont.)

- The ENDIT directory must reside on the same file system as the pool's data directory

- Several directories under the pool directory are recognized by the   ENDIT provider:

  ./request: ENDIT provider sends stage/migration requests here

  ./in: ENDIT provider checks for completed staged files in this directory

  ./out: ENDIT provider places files to be written to tape here

# ENDIT HPSSRetriever daemon

- HPSSRetriever is a daemon to stage filles from HPSS to dCache pools
- New software developed by BNL
- **Workflow**

  1. New stage requests are created under the ENDIT request directory by the ENDIT Provider.

  2. The daemon monitors the ENDIT request directory and detects incoming new stage requests.

  3. The daemon constructs stage requests and submits them to the ERADAT (HPSSBatch) request queue.

  4. ERADAT (HPSSBatch) sorts requests and stages files in batches.

  5. ERADAT (HPSSBatch) sends a callback after a file is in the HPSS disk cache.

  6. The daemon checks the callback content. If it's good, the daemon invokes PFTP to retrieve the file from HPSS.

  7. The file is staged from HPSS (disk cache) to the ENDIT ./in directory.

  8. ENDIT Provider detects the data file, moves it from ./in to the pool data directory, and marks the end of the stage request

**Brookhaven**
National Laboratory

# Staging Workflow with ENDIT HPSSRetriever



(1) Restore Request submission

(2) Watch the ENDIT request directory

(3) Construct stage requests and submit to HPSSBatch

(4) Sort requests and stage files in batch

(5) HPSSBatch sends callback after files in HPSS disk cache

(6) Issue PFTP get

(7) Stage file from HPSS (disk cache)

(8) Restore request done

**dCache Staging Pool**

/request

/in

**ENDIT HPSSRetriever Daemon**

**ERADAT (BNL HPSS Batch System)**

**HPSS**

**dCache ENDIT Provider Plugin**

**Tape System**

# Benefits of ENDIT HPSSRetriever

- Performance improvements on pool hosts

    - Eliminates polling on pool hosts, resulting in minimal load even with a high number of requests
    - Enables dCache to handle a large number of active staging requests simultaneously

- Provides flexible control over the maximum concurrent PFTP threads on each pool

- Prevents duplicated requests in HPSS Batch

- Reduces stress on PnfsManager and PoolManager due to non-polling nature

**Brookhaven** National Laboratory

# ENDIT HPSSArchiver Daemon

- HPSSArchiver is a daemon to flush dCache tape area files into HPSS.

- HPSSArchiver Daemon – Workflow

  1. New flush requests are created under the ENDIT request directory by the ENDIT Provider.

  2. The daemon monitors the ENDIT request directory and detects incoming new flush requests.

  3. The daemon invokes PFTP to flush files to HPSS.

  4. A file is flushed to HPSS successfully.

  5. ENDIT Provider detects the completion of a flushing process and marks the end of the flush request.

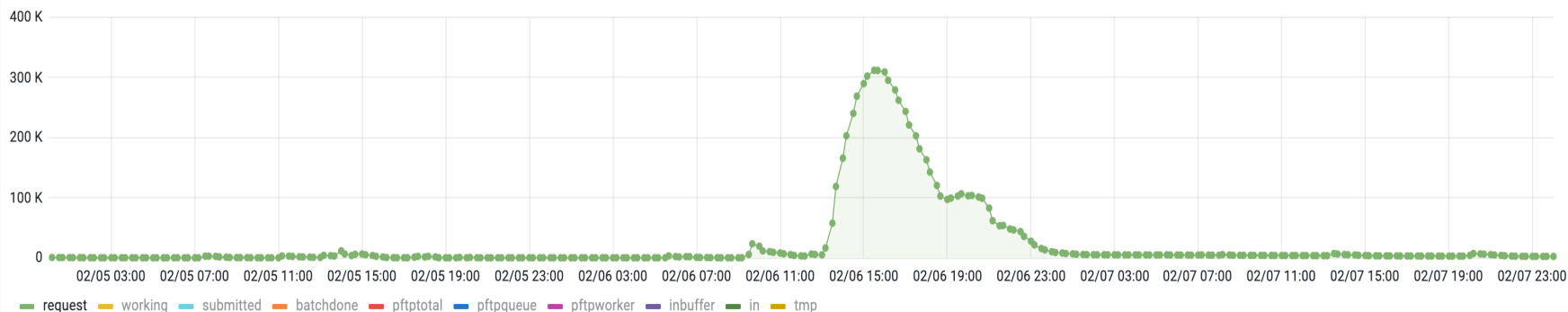Brookhaven™
National Laboratory

# Benefits of ENDIT HPSSArchiver

- Provides flexible control over the maximum concurrent PFTP threads on each pool

- A central place for a pool to control its migration requests. May add logic to do more (like hold requests for smart writing later ? ).

**Brookhaven** National Laboratory

# Restore Performance

# THANK YOU !

**Brookhaven** National Laboratory