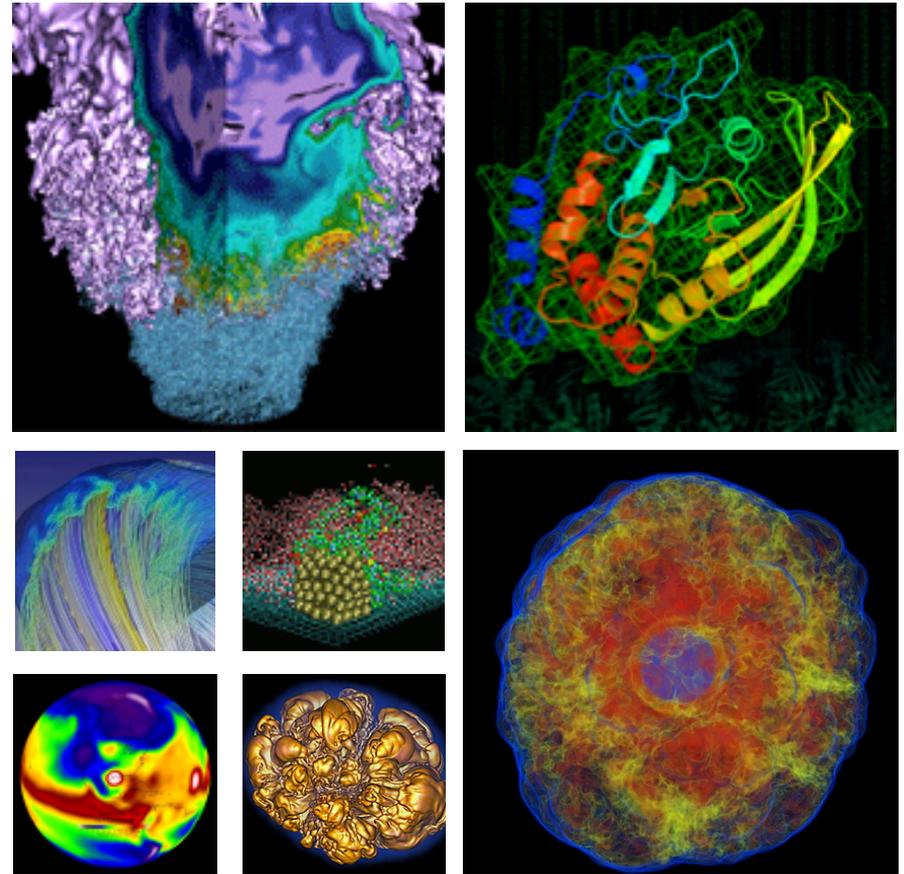


# Data Management at NERSC



**Lisa Gerhardt**  
NERSC User Services Group

**NUG Training**  
February 23, 2015

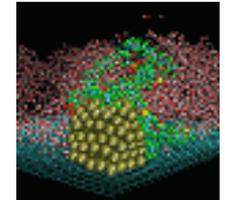
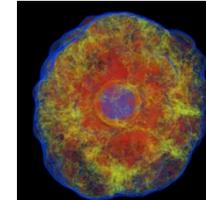
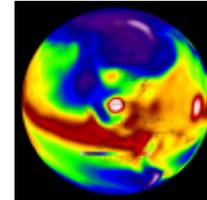
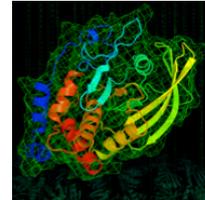
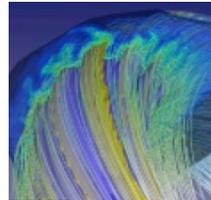
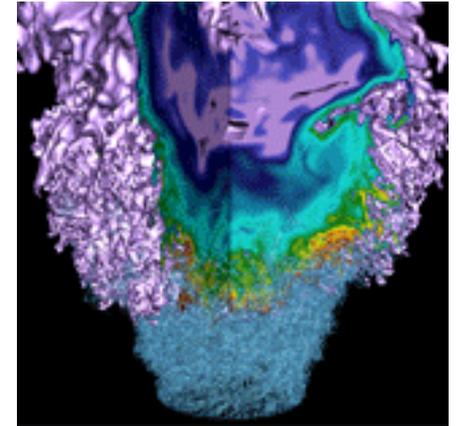
# Where Do I Put My Data?

---

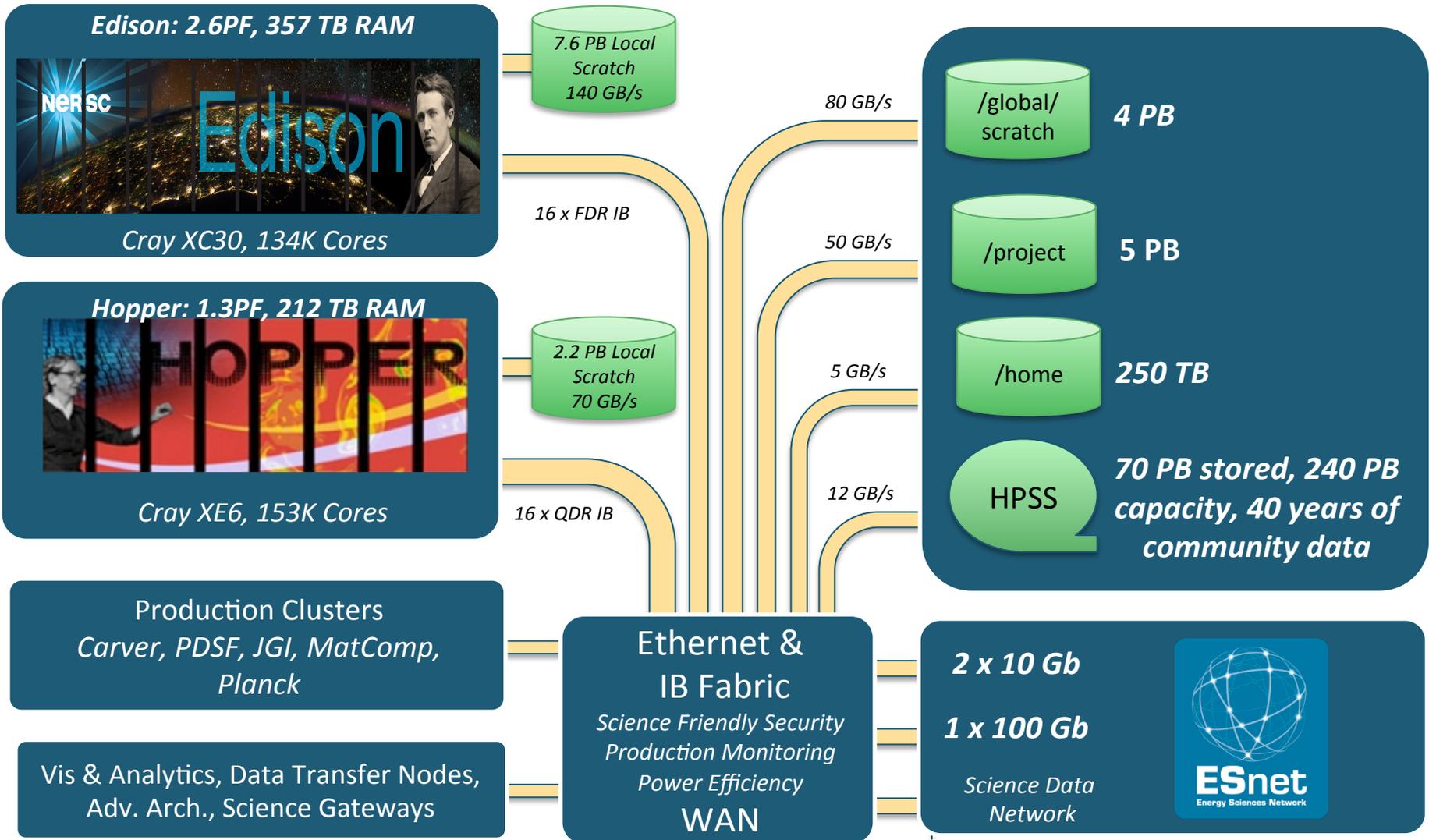


- **Overview of NERSC file systems**
  - Local vs. Global
  - Permanent vs. Purged
  - Personal vs. Shared
- **HPSS Archive System**
  - What is it and how to use it
- **Data Sharing**

# NERSC File Systems



# The compute and storage systems 2015



# Protect Your Data!

---



- Some file systems are backed up
- Some file systems are not backed up
- Restoration of individual files/directories may *not* be possible
- Hardware failures and human errors *will* happen

**BACK UP YOUR FILES TO HPSS!**

# Global File Systems



- **NERSC Global Filesystem (NGF)**
  - Based on IBM's General Parallel File System (GPFS)
- **Accessible on all NERSC systems**
- **Provides directories for home, global scratch, and project**
- **Shared by ~5000 active NERSC users**
  - Inefficient use affects others

# Global Homes File System Overview



- **Provided by two ~100 TB file systems**
  - 5 GB/s aggregate bandwidth
- **Access with \$HOME, ~/<file\_in\_home\_dir>**
- **Other name**
  - `/global/homes/l/lgerhard`
- **Low-level name**
  - `/global/u1/l/lgerhard`
  - `/global/u2/l/lgerhard -> /global/u1/l/lgerhard`

# Global Homes Use

---



- **Shared across all platforms**
  - Dot files that control user environment
  - `$HOME/edison`, `$HOME/hopper`, etc.
  - Visible only to you by default
- **Tuned for small file access**
  - Compiling/linking
  - Configuration files
  - **Do not send batch job output to \$HOME!**

- **Quotas enforced**
  - 40 GB
  - 1,000,000 inodes (i.e. files and directories)
  - Quota increases rarely (i.e., never) granted
  - Monitor with **myquota** command
- **“Permanent” storage**
  - No purging
  - Backed up (can retrieve files on your own with snapshots)
  - Hardware failures and human errors *will* happen

**BACK UP YOUR FILES TO HPSS!**

# Project File System Overview



- **Provides 5.1 PB high-performance disk**
  - 50 GB/s aggregate bandwidth
- **Available on all NERSC systems**
- **Intended for sharing data between platforms, users, or with the outside world**
- **Beginning this year every MPP repo gets a project directory**

`/project/projectdirs/m9999`

- **Tuned for large streaming file access**
  - Sharing data within a project or externally
  - Running I/O intensive batch jobs
  - Data analysis/visualization
- **Access controlled by Unix file groups**
  - Visible, writable by whole group by default
  - Group name usually same as directory
  - Requires administrator (usually the PI or PI Proxy)
  - Can also use access control list (ACL)

- **Quotas enforced**
  - 1 TB
  - 1,000,000 inodes
  - Quota increases may be requested
  - Monitor with `prjquota` command
    - `% prjquota bigsci`
- ***Permanent storage***
  - No purging
  - Backed up if quota  $\leq$  5 TB
  - Hardware failures and human errors *will* happen

**BACK UP YOUR FILES TO HPSS!**

# Global Scratch File System Overview



- **Provides 4 PB high-performance disk**
  - 80 GB/s aggregate bandwidth
- **Access with \$GSCRATCH**
- **Low-level name**  
`/global/scratch2/sd/lgerhard`

# Global Scratch Use

---



- **Shared across all systems**
  - Visible only by you by default
  - Primary scratch file system for Carver
- **Tuned for large streaming file access**
  - Running IO intensive batch jobs
  - Data analysis/visualization

# Global Scratch Policies



- **Quotas enforced**
  - 20 TB
  - 4,000,000 inodes
  - Quota increases may be requested
  - Monitor with `myquota` command
- **Temporary storage**
  - Bi-weekly purges of *all* files that have not been accessed in over 12 weeks
    - List of purged files in `$GSCRATCH/.purged.<timestamp>`
  - Hardware failures and human errors *will* happen

**BACK UP YOUR FILES TO HPSS!**

# Local File Systems on Cray Machines



- **Edison and Hopper have local scratch**
- **Edison has two *scratch* file systems**
  - Users randomly assigned
  - Each has 2.1 PB (1 PB on Hopper)
  - Each has 48 GB/s aggregate bandwidth (35 GB/s Hopper)
- **Edison has extra high-performance scratch (scratch3)**
  - 3.2 PB, 72 GB/s aggregate bandwidth
- **Provided by Cray, based on Lustre**
- **Generally, IO access for batch jobs on Hopper and Edison will be fastest for local scratch**

# Edison Scratch Use



- Each user gets a scratch directory in `/scratch1` or `/scratch2` (Hopper: `/scratch` or `/scratch2`)  
`/scratch2/scratchdirs/dptturner`
  - Best name: `$SCRATCH`
  - Visible only to you by default
- Access to `/scratch3` must be requested
  - Large datasets
  - High bandwidth
- Tuned for large streaming file access
  - Running I/O intensive batch jobs
  - Data analysis/visualization

# Edison Scratch Policies



- **Quotas enforced in \$SCRATCH by submit filter**
  - 10 TB (5 TB Hopper)
  - 10,000,000 inodes (5M inodes Hopper)
  - Quota increases may be requested
  - Monitor with **myquota** command
  - No quota enforcement in /scratch3
- **Temporary storage**
  - Daily purges of *all* files that have not been accessed in over 12 weeks (8 weeks on scratch3)
    - List of purged files in \$SCRATCH/.purged.<timestamp>
  - Hardware failures and human errors *will* happen

**BACK UP YOUR FILES TO HPSS!**

- **Global home directories (\$HOME)**
  - Source/object/executable files, batch scripts, input files, configuration files, batch job summaries (*not* for running jobs)
  - Backed up
  - 40 GB permanent quota
- **Global project directories**
  - Sharing data between people and/or systems, short term data storage
  - Backed up if quota less than or equal to 5 TB
  - All MPP repos have one, 1 TB default quota

# Short-Term File Systems



- **Local scratch directories**
  - Cray (Edison, Hopper) only
  - Large, high-performance parallel Lustre file system
  - Not backed up; files purged after 12 weeks
  - Hopper: 5 TB default quota; Edison: 10 TB default quota
  - `$SCRATCH`, `$SCRATCH2`
- **Global scratch directories**
  - All systems
  - Large, high-performance parallel GPFS file system
  - Not backed up; files purged after 12 weeks
  - 20 TB default quota
  - `$GSCRATCH`

# Where Do I Put My Data?

## Local Scratch

Fastest IO

Only visible on one machine

Only visible to you

Purged



## Project

Medium IO

Visible on all machines

Visible to all group members

Never purged

External sharing

## Global Scratch

Fast IO

Visible on all machines

Only visible to you

Purged

## Home

Source code, config. files

Only visible to you

**No batch output**

# File Systems Summary



File System	Path	Type	Default Quota	Backups	Purge Policy
Global Homes	\$HOME	GPFS	40 GB / 1M inodes	Yes	Not purged
Global Scratch	\$GSCRATCH	GPFS	20 TB / 4M inodes	No	12 weeks from last access
Global Project	/project/ projectdirs/ projectname	GPFS	1 TB / 1M inodes	Yes, if quota less than or equal to 5TB	Not purged
Hopper Scratch	\$SCRATCH and \$SCRATCH2	Lustre	5 TB / 5M inodes (combined)	No	12 weeks from last access
Edison Scratch	\$SCRATCH	Lustre	10 TB / 5M inodes (none in /scratch3)	No	12 weeks from last access

# Resources

---



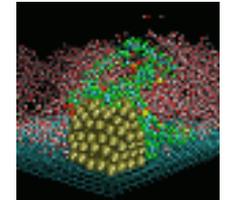
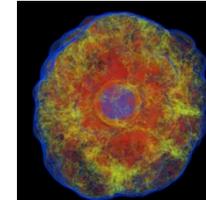
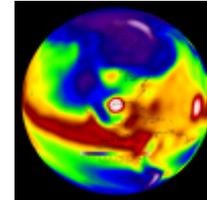
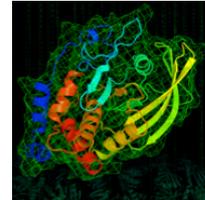
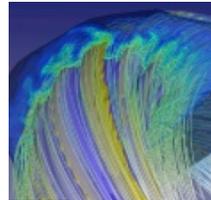
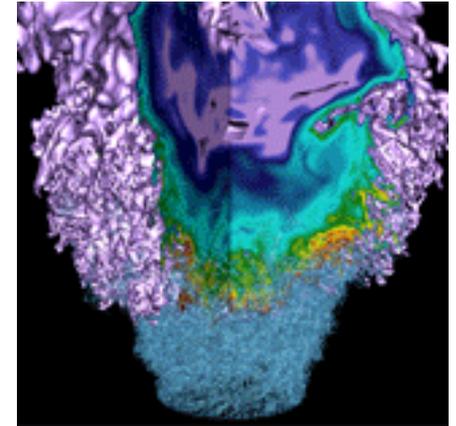
<http://www.nersc.gov/users/data-and-file-systems/>

[http://www.nersc.gov/users/data-and-file-systems/  
file-systems/](http://www.nersc.gov/users/data-and-file-systems/file-systems/)

[http://www.nersc.gov/users/computational-systems/  
edison/file-storage-and-i-o/](http://www.nersc.gov/users/computational-systems/edison/file-storage-and-i-o/)

[http://www.nersc.gov/users/computational-systems/  
hopper/file-storage-and-i-o/](http://www.nersc.gov/users/computational-systems/hopper/file-storage-and-i-o/)

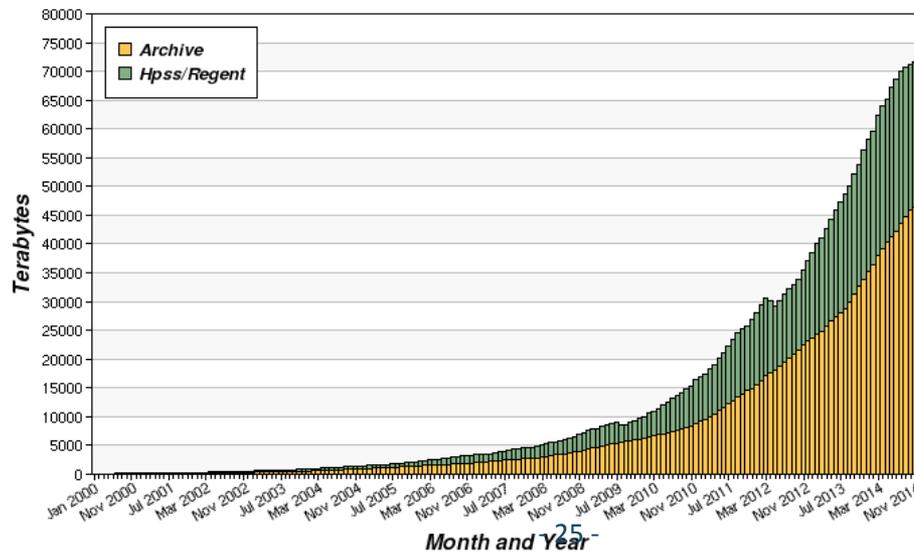
# HPSS: The NERSC Archive System



# Archiving Data is Necessary

- **Data growth is exponential and file system space is finite**
  - 80% of stored data is never accessed after 90 days
  - Cost of storing infrequently accessed data on flash or spinning disk is prohibitive
  - Store important data in an archive to free faster resources for processing workload
    - Data from publications, unique experimental, or simulation data
- **NERSC provides the HPSS archive system for data archiving**

*Cumulative Storage by Month and System*



**70 PB of data**  
Started in 1998, but  
oldest file is from  
the 70s

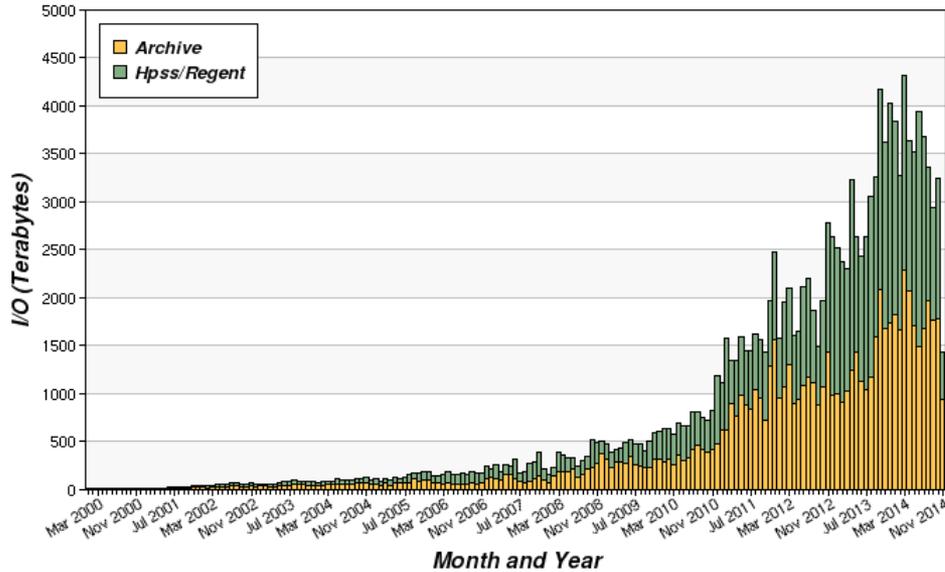
# Features of the NERSC archive



- **NERSC implements an online or “active archive”**
  - Parallel high-speed transfer and fast data access
    - Data is transferred over parallel connections to the NERSC internal 10Gb network
    - Access to first byte in seconds or minutes as opposed to hours or days
  - Tiered internal storage facilitates high speed data access:
    - Initial data ingest to high-performance disk cache
    - Data migrated to automated enterprise tape system and managed by HSM software (HPSS) based on file age and usage
  - Indefinite data retention policy
- **The archive is accessible to all NERSC users**
- **Often referred to as HPSS**

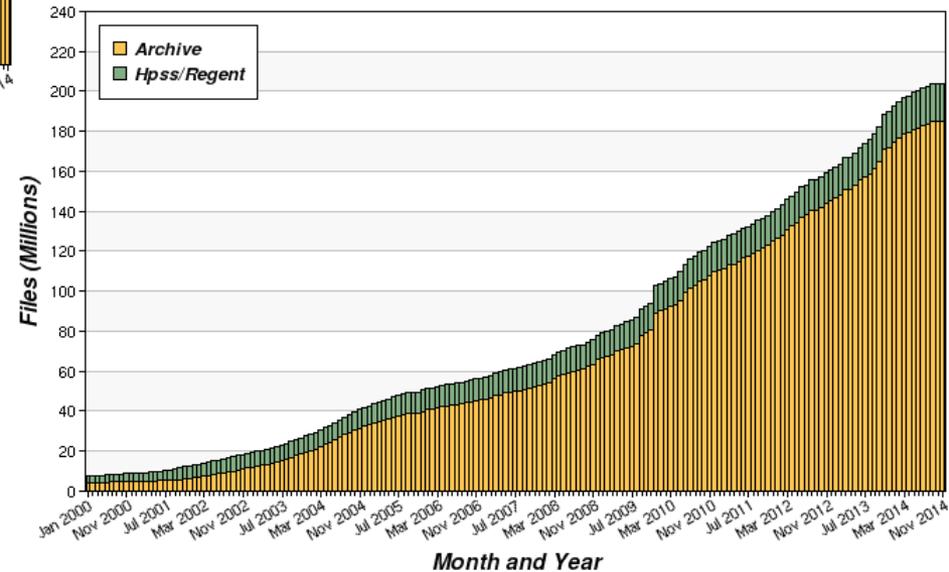
# HPSS is Heavily Used

Monthly I/O by Month and System



~100 TB / day

Number of Files Stored by Month and System



# Accessing HPSS from NERSC Systems



- **HSI**

- Fast, parallel transfers, unix-like interface
- Store from file system to archive:

```
-bash-3.2$ hsi
```

```
A:/home/n/nickb-> put myfile
```

```
put 'myfile' : '/home/n/nickb/myfile' ( 2097152 bytes, 31445.8 KBS (cos=4))
```

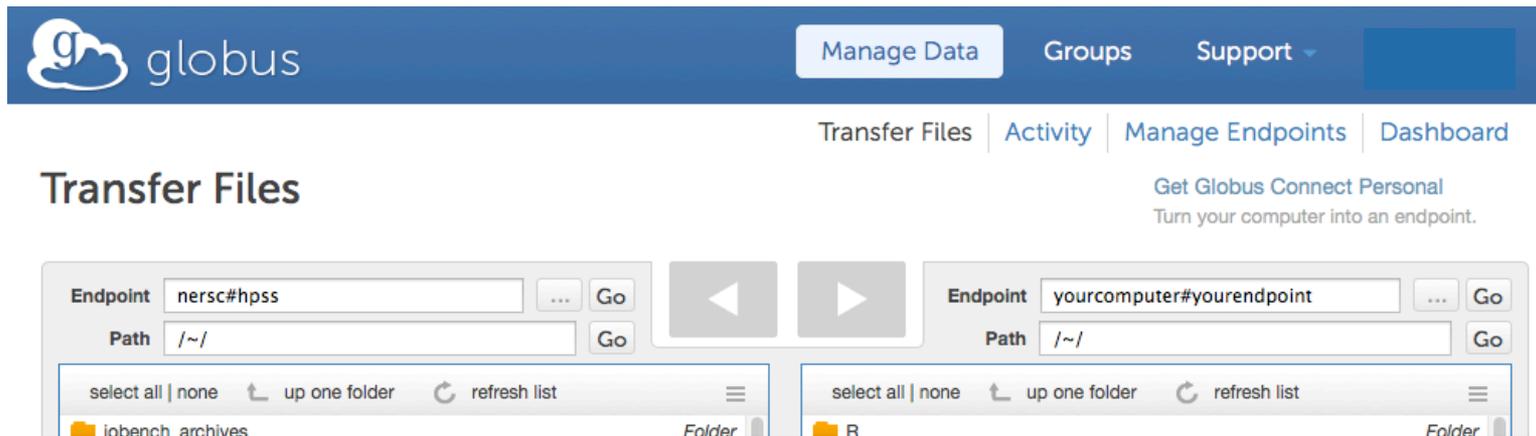
- **HTAR**

- Parallel, puts tar file directly into HPSS, excellent for groups of small files
- Syntax: *htar [options] <archive file> <local file | dir>*

```
-bash-3.2$ htar -cvf /home/n/nickb/mydir.tar ./mydir
```

# Accessing HPSS from Outside NERSC

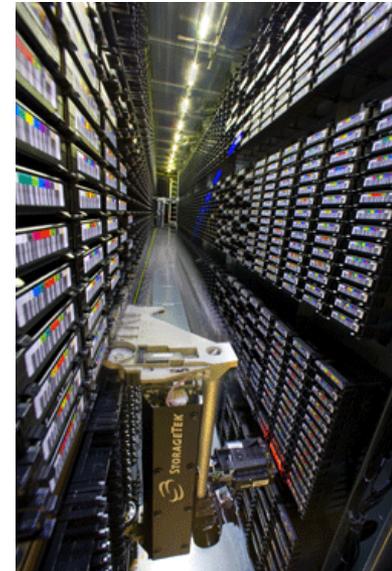
- HSI and HTAR precompiled binaries available for most systems
- **ftp**: non-parallel, but common
- **gridFTP**: parallel, requires credential
- **Globus**: parallel, requires endpoint



The screenshot shows the Globus Connect Personal web interface. At the top, there is a navigation bar with the Globus logo, a search bar, and links for 'Manage Data', 'Groups', and 'Support'. Below this, there are links for 'Transfer Files', 'Activity', 'Manage Endpoints', and 'Dashboard'. The main content area is titled 'Transfer Files' and includes a link to 'Get Globus Connect Personal' with the text 'Turn your computer into an endpoint.' Below this, there are two side-by-side endpoint configuration panels. Each panel has an 'Endpoint' field, a 'Path' field, and 'Go' buttons. The left panel shows 'Endpoint: nersc#hpss' and 'Path: /~/'. The right panel shows 'Endpoint: yourcomputer#youendpoint' and 'Path: /~/'. Below the endpoint fields, there are controls for file selection: 'select all | none', 'up one folder', and 'refresh list'. The left panel shows a folder named 'lobench archives' and the right panel shows a folder named 'R'.

# Tape IO Characteristics

- **Ultimately all data in HPSS is written to tape**
- **Tape is linear media**
  - Behaves differently than disk:
    - Data cannot be re-written in place, it is appended at the end
    - Reading and writing are sequential operations – no random access
- **Tape drives behave differently than disk drives**
  - Take time to seek to file locations on tape
  - Take time to ramp up to full speed
  - Tape drives stop after reading or writing each file
- **HPSS will not respond like a normal file system**
  - Presents itself as one, but some things can have unexpected results



# Size Matters

- **Sweet Spot**
  - Tape drives perform best when operating at full rate for long durations
  - Large file are best for tape drive performance
  - Many small files causes frequent stops and low-speed operations, can take a very long time to retrieve
  - File bundles in the **100s of GB** currently provide best performance
- **Group small files for optimal storage**
  - Use HTAR, GNU tar, or zip to bundle groups of small files together to optimize tape and network performance
- **There is such a thing as too big**
  - Files spanning multiple tapes incur tape mount delays

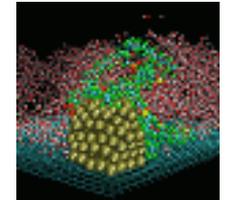
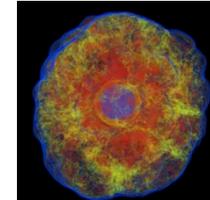
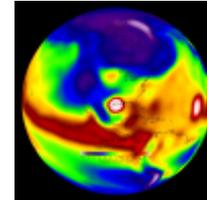
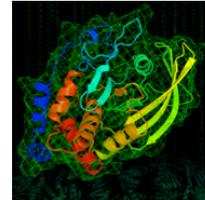
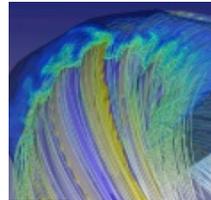
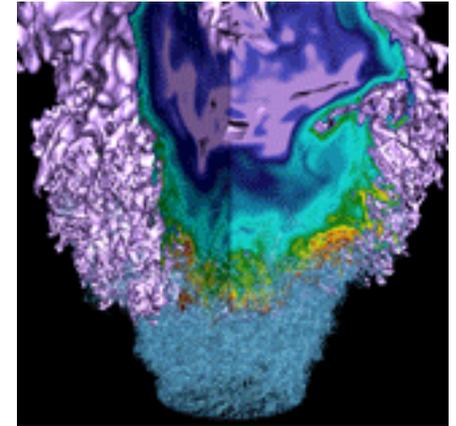
- **Group small files together and avoid excessive writes**
  - Use htar or tar to group into ~100s of GB
- **Order your retrievals**
  - Grab files from a tape in order of tape position
  - Grab all files from a tape while tape is mounted
- **Avoid excessive transfer failures**
  - Globus with unreliable network will retry many times
  - Directory permission issues
- **No exclusive access to the archive**
  - No batch system
  - Inefficient use affects performance for everyone

# Further Reading



- **NERSC Website**
  - Archive documentation:
    - <http://www.nersc.gov/users/data-and-file-systems/hpss/getting-started/>
  - Data management strategy and policies:
    - <http://www.nersc.gov/users/data-and-file-systems/policies/>
  - Accessing HPSS
    - <http://www.nersc.gov/users/data-and-file-systems/hpss/getting-started/accessing-hpss/>
- **HSI and HTAR man pages are installed on NERSC compute platforms**
- **Gleicher Enterprises Online Documentation (HSI, HTAR)**
  - <http://www.mgleicher.us/index.html/hsi/>
  - <http://www.mgleicher.us/index.html/htar/>
- **“HSI Best Practices for NERSC Users,” LBNL Report #LBNL-4745E**
  - [http://www.nersc.gov/assets/pubs\\_presos/HSIBestPractices-Balthaser-Hazen-2011-06-09.pdf](http://www.nersc.gov/assets/pubs_presos/HSIBestPractices-Balthaser-Hazen-2011-06-09.pdf)

# Data Sharing



# Data Sharing

---

- **Ensure security**
  - Do not share passwords
  - Do not share files from \$HOME
- **Project directories designed for sharing**
  - Open to anyone in the repository
- **Use Unix *group* permissions**
  - Request creation of Unix group
  - Set permissions with chgrp/chmod
    - Use setgid bit

# give/take

- New, but based on *old* LLNL and LANL commands
- Appropriate for smaller files

```
joe% give -u bob coolfile
```

- File copied *to* spool location
- Bob gets email telling him Joe has given him a file

```
bob% take -u joe coolfile
```

- File copied *from* spool location

- Spooled files count against *giver's* GSCRATCH quota

# Science Gateways on Project



- **Make data available to outside world**

```
mkdir /project/projectdirs/bigsci/www
```

```
chmod o+x /project/projectdirs/bigsci
```

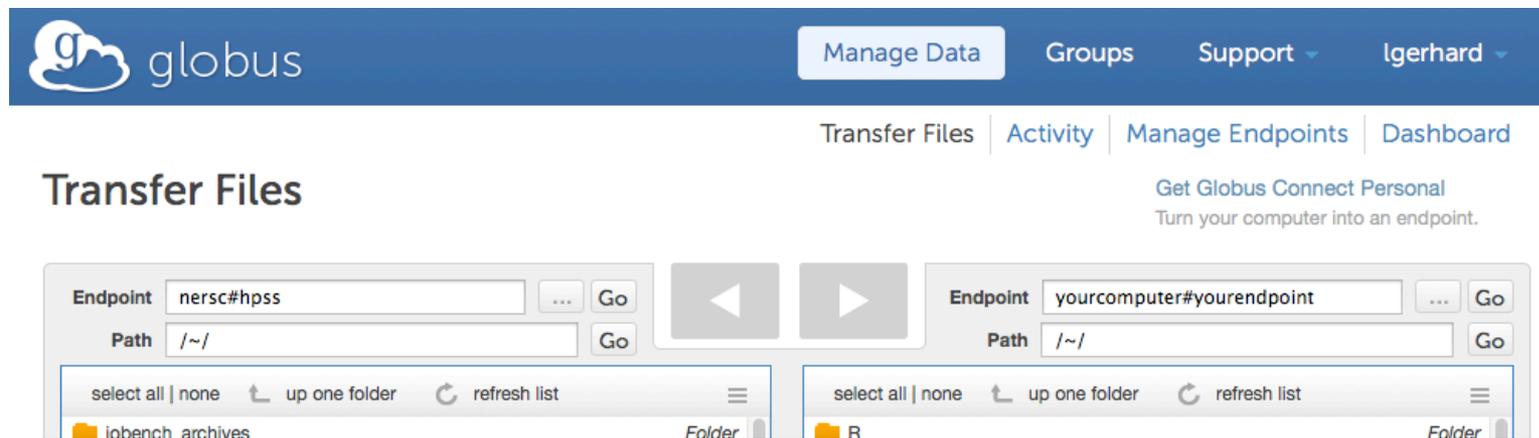
```
chmod o+rx /project/projectdirs/bigsci/www
```

- **Access with web browser**

```
http://portal.nersc.gov/project/bigsci
```

- **Global file systems**
    - Use *local* cp instead of *remote* scp
  - **Use scp for small-to-medium files over short-to-medium distance**
    - Even better if HPN versions installed
- ```
% ssh -v  
OpenSSH_5.1p1NMOD_2.9-hpn13v5, OpenSSL 0.9.8e-fips-rhel5 01 Jul  
2008
```
- **Use bbcp for larger files and/or longer distances**
    - Many tuning options
    - Complicated command line

- **Do-it-all web-based file transfer service**
- **High-performance**
  - Parallel data channels (gridftp)
- **Fire and forget model**
- **Also has a command-line interface for scripting**



globus Manage Data Groups Support lgerhard

Transfer Files Activity Manage Endpoints Dashboard

Transfer Files [Get Globus Connect Personal](#)  
Turn your computer into an endpoint.

Endpoint: nersc#hpss Go Path: /~/ Go

Endpoint: yourcomputer#yourendpoint Go Path: /~/ Go

select all | none up one folder refresh list Folder

select all | none up one folder refresh list Folder

# Further Reading

---

- **Sharing data**

- <https://www.nersc.gov/users/data-and-file-systems/sharing-data/>

- **Transferring Data**

- <https://www.nersc.gov/users/data-and-file-systems/transferring-data/>

# Asking Questions, Reporting Problems

---



- **Contact NERSC Consulting**
  - Toll-free 800-666-3772
  - 510-486-8611, #3
  - Email [consult@nersc.gov](mailto:consult@nersc.gov).
  - <https://www.nersc.gov/users/getting-help/>



**Thank you.**