# NCAR Storage Spaces

### And managing your data holdings

**Mick Coady**,
*CISL Consulting Services*
mickc@ucar.edu

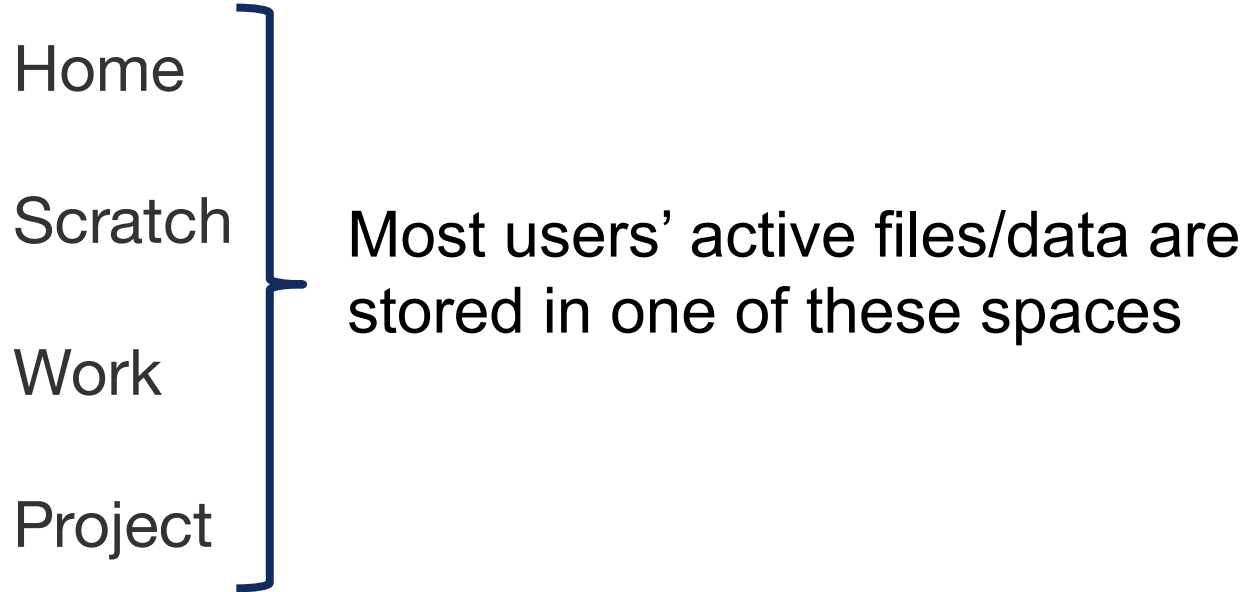**April 22, 2020**

# What we'll cover

- GLADE file spaces
  - home
  - scratch
  - work
  - project

- Campaign Storage

- HPSS
  - historical perspective
  - migration urgency

- Collections

- Managing your data

# GLADE

- GLADE - **Gl**obally **A**ccessible **D**ata **E**nvironment

- High-performance shared file system

- Available across all of CISL's HPC systems
  - Cheyenne
  - Casper
  - data-access nodes

# GLADE spaces

Home

Scratch

Work

Project

Most users' active files/data are stored in one of these spaces

# GLADE home

Home

Scratch

Work

Project

- Every user has a home directory
  - /glade/u/home/*<userid>*

- 25 GB of <u>useable</u> storage
  - *Don't be fooled by gladequota!*

- Generally used for scripts, source code, and small data sets

# GLADE home (cont.)

Home

Scratch

Work

Project

- Backups
  - Backed up several times/week
  - Each backup is retained for several weeks
  - Contact CISL to recover a backup

- Snapshots
  - Created several times/day (if there are changes)
  - Run the command **snapls** for a list of your available snapshots
  - For more information see: https://www2.cisl.ucar.edu/resources/storage-and-file-systems/glade/recovering-files-snapshots

# GLADE scratch

Home

**Scratch**

Work

Project

- Every user has a scratch directory
  - /glade/scratch/*<userid>*

- 15 PB total capacity

- 10 TB default user quota
  - users can request larger, temporary extensions
  - justifications required
  - CISL reserves right to reduce extensions

# GLADE scratch (cont.)

Home

Scratch

Work

Project

- Temporary storage space

- Good for large output files

- 120-day purge policy
  - Based on file's last access time
  - Subject to change

- NOT backed up!
- NOT backed up!
- NOT backed up!

# GLADE work space

Home

Scratch

Work

Project

- Every user has a work directory
  - /glade/work/*<userid>*

- 1 TB quota

- Good for active files needing longer retention than scratch space, compiled executables, etc.

- Not purged

- NOT backed up!
- NOT backed up!
- NOT backed up!

# GLADE project space

Home

Scratch

Work

Project

- Each NCAR lab has a high-level directory
  e.g.  /glade/p/cisl   /glade/p/cgd   /glade/p/mmm
         /glade/p/hao   /glade/p/eol   /glade/p/acom

- Many university projects have a high-level directory
  e.g.  /glade/p/univ/*<projectname>*
         /glade/p/uwyo/*<projectname>*

- NCAR labs and universities determine their organization's directory sub-structure

# GLADE project space (cont.)

Home

Scratch

Work

Project

- No default quotas

- CISL User Services manages allocations

- NCAR labs and universities are responsible for managing their usage

- Not purged

- NOT backed up!
- NOT backed up!
- NOT backed up!

# Campaign Storage

- Cost-effective, medium-performance disk-based file system

- Warm archive for storing data on publication timescales
  - ~ 5 years

- Available as a Globus endpoint
  - "NCAR Campaign Storage"

- Available from Casper and data-access
  - Mounted as /glade/campaign

- Not available from Cheyenne
  - To use data on Cheyenne, copy to a mounted file system
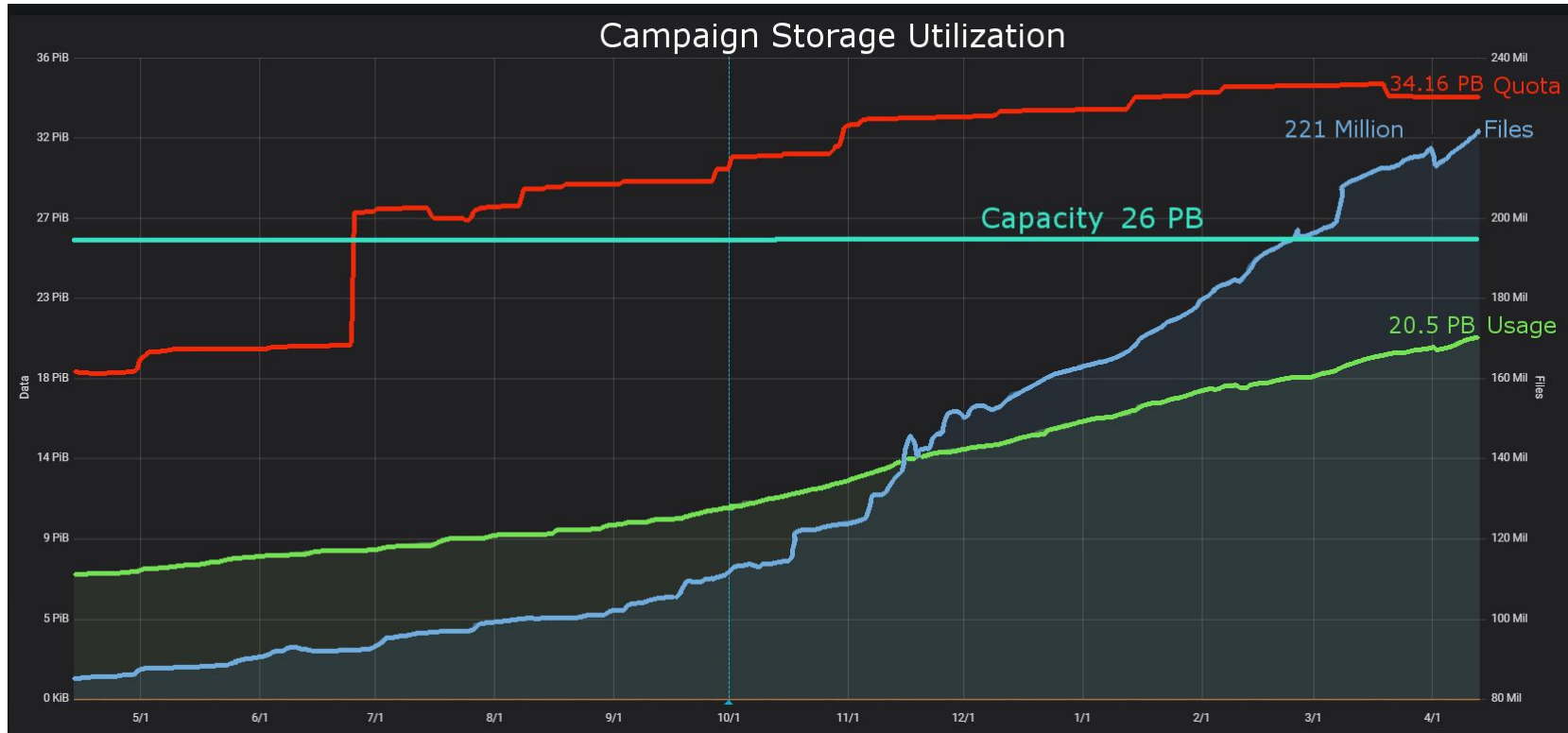
# Campaign Storage (cont.)

- 43 PB total capacity  (as of 4/20/20)

- Currently over 21 PB of data stored
  - 48% full

- More than 365 users

- Planned annual augmentation, pending budget <u>and</u> utilization

- Second augmentation late summer 2020
  - NSF approval received 4/21/20
  - An additional 24 PB

# Campaign Storage (cont.)

- Each NCAR Lab and many university projects have a high-level directory
  - Similar organization as GLADE Project space

- CISL User Services manages allocations
  - No default quotas

- NCAR Labs and universities are responsible for managing their usage
  - Labs and universities can purchase private Campaign Storage space
    - Contact CISL User Services for pricing and process details

- 5 year purge policy
  - Based on a file's <u>creation</u> date

- Not backed up!   Not backed up!   Not backed up!

# Campaign Storage (cont.)



Campaign Storage Utilization

34.16 PB Quota

221 Million Files

Capacity 26 PB

20.5 PB Usage

NCAR UCAR | NCAR Storage Spaces

# HPSS

- Tape archive file system

- Cold archive for long-term data storage

- Now in end-of-life
  - Will be decommissioned on Oct 1, 2021
  - Technically and financially no longer unsustainable
  - Became read-only on Jan 20, 2020

- Replaced by Campaign Storage

- CISL is actively working with users to vacate
  - Started with more than 90 PB of data and over 300 million files

# HPSS (cont.)

- Will be decommissioned on Oct 1, 2021
- Still more than 17 months away
- So there's plenty of time to migrate my data - right?

- ## No!

- ~80 PB of data remains on HPSS
- At an aggressive 2 PB/month read rate
- Only 42% of data can still be read

- If 80 PB still remains on 10/1/20 ...
  - only 30% will be able to be read
- If 80 PB still remains on 4/1/20 ...
  - only 15% will be able to be read

# Collections

- Reserved for a small number of permanent collections

- Mounted as /glade/collections

- Curated research data that includes:
  - Research Data Archive (RDA)
  - Climate Data Gateway (CDG)
  - CMIP AP
  - Earth Observatory Lab (EOL) field campaigns
  - High Altitude Observatory (HAO)

- Available by special request only
  - CISL User Services manages requests

# Managing your data

- gladequota command

- access_report.txt files
  - Campaign Storage
  - Project space directories

- GUFI

Tools readily available to all users to help manage their active data holdings

DASH

For tomorrow and beyond

# Managing your data
## gladequota command

- Quick and easy to use
  - Maintained by CISL
  - GLADE spaces updated real-time
  - Campaign Storage updated hourly

- Returns relevant information on storage spaces you have write access to

- If you reach your disk quota on any of your GLADE files spaces you may experience weird problems
  - Can't log on to systems
  - Unexplained job failures
  - Particularly true for home and scratch directories
  - Run gladequota command regularly to check your usage

# Managing your data
## gladequota command example

```
mickc@cheyenne1:~-> gladequota
Current GLADE space usage: mickc
```

Beware!

| Space | Used | Quota | % Full | # Files |
|-------|------|-------|--------|---------|
| /glade/scratch/mickc | 143.36 GB | 10.00 TB | 1.40 % | 2308 |
| /glade/work/mickc | 273.34 GB | 1024.00 GB | 26.69 % | 63763 |
| /glade/u/home/mickc | 12.72 GB | 50.00 GB | 25.44 % | 66993 |
| /glade/u/benchmarks | 94.30 GB | 250.00 GB | 37.72 % | 55 |
| /glade/u/apps | 762.91 GB | 1024.00 GB | 74.50 % | 7731127 |
| /glade/p/cisl/CSG | 1.69 TB | 5.00 TB | 33.79 % | 8168998 |
| Campaign: mickc (user total) | 68.85 GB | n/a | n/a | |
| /glade/campaign/cisl/csg | 4.19 TB | 23.00 TB | 18.22 % | 1585 |

```
(Campaign usage as of: Wed Apr  8 11:05:04 MDT 2020)

/glade/scratch  - 70.4% used (10564 TB used out of 15000 TB total)
```

NCAR
UCAR

# Managing your data
## access_report.txt files

- Available for /glade/p and Campaign Storage

- Provides a quick look at data access trends

- Located in directories under the lab and university levels, e.g.
  - /glade/p/cisl/CSG/access_report.txt
  - /glade/p/uwyo/wyom0085/access_report.txt
  - /glade/campaign/cgd/access_report.txt

- Frequently used by lab and university data managers
  - but accessible by everyone with read access

# Managing your data
## access_report.txt files - example

```
mickc@cheyenne6:-> cat /glade/p/cisl/CSG/access_report.txt
Directory: /glade/p/cisl/CSG
Total Files: 8.17 M
Total Data: 1.68 TB
Scan date: 2020-03-30


Last Accessed          Data (%)              # Files (%)
-------------------------------------------------------------
<1 Month             2.32 GB (0.13%)         845,663 (10.35%)
1 Month                 0 b (0.00%)               0 (0.00%)
6 Months            62.01 MB (0.00%)          11,427 (0.14%)
1 Year               1.67 TB (99.86%)          7.31 M (89.51%)
3 Years                 0 b (0.00%)               0 (0.00%)
5+ Years                0 b (0.00%)               0 (0.00%)

  Owner                Data (%)              # Files (%)
-------------------------------------------------------------
<1 Month:
  1. csgteam         2.26 GB (0.13%)         719,031 (8.80%)
  2. ddvento        61.33 MB (0.00%)         126,616 (1.55%)
  3. valent         99.02 KB (0.00%)              11 (0.00%)
  4. vanderwb        8.00 KB (0.00%)               2 (0.00%)
  5. sghosh          4.00 KB (0.00%)               1 (0.00%)
```

# Managing your data
## access_report.txt files – example (cont.)

```
1 Month:
6 Months:
  1. csgteam       48.37 MB (0.00%)        11,418 (0.14%)
  2. ddvento       13.63 MB (0.00%)             1 (0.00%)
  3. vanderwb       8.00 KB (0.00%)             7 (0.00%)
  4. valent          376 b (0.00%)             1 (0.00%)
1 Year:
  1. ddvento        1.13 TB (67.56%)       465,932 (5.70%)
  2. csgteam      554.56 GB (32.30%)        6.85 M (83.79%)
  3. valent        11.29 MB (0.00%)           945 (0.01%)
  4. bjsmith      221.92 KB (0.00%)             2 (0.00%)
3 Years:
5+ Years:
```

# Managing your data
## GUFI

- **G**rand **U**nified **F**ile **I**ndex

- New tool deployed by CISL to help users manage their data holdings

- Developed at Los Alamos National Lab

- Provides extremely fast queries of massive file systems

- Documentation: https://github.com/mar-file-system/GUFI

# Managing your data
## GUFI (cont.)

- CISL updates the GUFI database weekly

  - ○ GLADE – home, scratch, project, and work

  - ○ Campaign Storage

  - ○ HPSS
    - ■ CISL uses GUFI to generate weekly HPSS user and project file lists

- Users can write their own SQL-like queries but …

- CISL provides easy-to-use query wrapper scripts

# Managing your data
## GUFI wrapper scripts

- Simplifies many common user queries
- Available on Casper (not Cheyenne)
- Currently being evaluated by NCAR Storage Advisors Group
- Two step process:
  - **qdh**
    - Generates a query to GUFI server
    - Stores output files in user's scratch directory

  - *qdh_plot*
    - *under construction - coming soon*

# Managing your data
## GUFI (cont.)

To run CISL's GUFI query wrapper scripts

1.  log on to Casper (**casper.ucar.edu**)
2.  load the gufiwrappers module:  *module load gufiwrappers*
3.  run command *qdh* with appropriate parameters and options
4.  *When it becomes available, run **qdh_plot** to produce graphical reports*

Documentation available on github

https://ncar.github.io/gufiwrappers/index.html

# Managing your data
## Simple GUFI example

```
mickc@casper26:~-> module load gufiwrappers
mickc@casper26:~-> qdh --storage=campaign --directory=/glade/campaign/cisl/csg --list=filename,size,owner --filter-by-users=mickc
-------------------------------------------------------------------------------------
Writing log file... /gpfs/fs1/scratch/mickc/gufi_tmp/logs/log_20200421_083415.log
Executing gufi command...writing cache files in:
   /gpfs/fs1/scratch/mickc/gufi_tmp/raw/campaign/cisl/__.csg.dat.*          <──── Output is not intended for humans
--------------------------------------------------------------------------------------------
--------------------------------------------------------------------------------------------
Writing output in file:
/gpfs/fs1/scratch/mickc/gufi_tmp/reports/rep_20200421_083415.dat
--------------------------------------------------------------------------------------------
Connection to squall.ucar.edu closed.
```

# Managing your data
## Simple GUFI example (cont.)

```
mickc@casper26:  cd /glade/scratch/mickc/gufi_tmp/reports
mickc@casper26:  /glade/scratch/mickc/gufi_tmp/reports-> cat rep_20200421_083415.dat

/glade/campaign/cisl/csg/mickc/pnichols-glade-u-home-backup-2018-12-04.tar.gz,13079279358,mickc
/glade/campaign/cisl/csg/mickc/pnichols_old_work.tar.gz,60759557809,mickc
/glade/campaign/cisl/csg/mickc/python_scripts.tar,95427,mickc
/glade/campaign/cisl/csg/mickc/queue_status_ch.tar.gz,4235,mickc
/glade/campaign/cisl/csg/mickc/wrfbdy_d01,90392564,mickc
```

Remember - qdh_plot is on the way! 😉

# DASH and DSET

- **D**igital **A**sset **S**ervices **H**ub
- **D**ata **S**tewardship **E**ngineering **T**eam

- UCAR Policy 3-5 "Publication & Information Dissemination" - Oct 2009

  "UCAR supports an open exchange of data and scholarly information derived from our research. It is UCAR's policy to share this scientific and technical information with the community…"

- UCAR Procedure to Support Open Access Data Requirements - May 2019

  "NCAR and UCP project PIs **must** engage the Digital Asset Services Hub (DASH) during the initial planning/budgeting phase of research projects to: 1) determine if data management services will be needed to support the project, and 2) if needed, select a data management service tier and develop a data management services budget to be included in the project proposal."

# DASH and DSET (cont.)

Project planning to meet Open Access Data Requirements

1. Determine if data management services will be needed to support the project
   - Work through the "Determine Data Management Requirements for Proposals"
   - See https://www2.cisl.ucar.edu/dash/data-management-guidance

2. If data management services are needed to support the project, submit a request for repository support at https://submit-data.ucar.edu

   NCAR Dataset Submission Request System
   https://submit-data.ucar.edu/

# DASH and DSET (cont.)

Bottom line – plan ahead!

CISL's DASH/DSET team:
- Doug Schuster    - Manager, CISL DECS
- Bridget Thrasher - DASH Coordinator
- Matt Mayernik    - DSET Chair

# NCAR Storage Spaces

# Questions?