# RESEARCH DATA DEPOT: INFRASTRUCTURE FOR DATA

**PURDUE** UNIVERSITY

**Preston Smith**

Manager of Research Computing Support

# TODAY'S AGENDA

- Welcome  (Donna Cumberland, Executive Director, Research Computing)

- Introduction (Dr. Gerry McCartney, System CIO)

- Research Data Depot (Preston Smith, Manager of Research Computing Support)

- 2014-2015 Computation Plans (Michael Shuey, HPC Technical Architect)

**PURDUE**
UNIVERSITY

**Since Steele in 2008, Research Computing has deployed many world-class offerings in computation**

PURDUE
UNIVERSITY

# SIX COMMUNITY CLUSTERS

## STEELE

7,216 cores

Installed May 2008

Retired Nov. 2013

## COATES

8,032 cores

Installed July 2009

24 departments

61 faculty

Retired Sep. 2014

## ROSSMANN

11,088 cores

Installed Sept. 2010

17 departments

37 faculty

## HANSEN

9,120 cores

Installed Sept. 2011

13 departments

26 faculty

## CARTER

10,368 cores

Installed April 2012

26 departments

60 faculty

#175 on June 2013 Top 500

## CONTE

9,280 Xeon cores
(69,600 Xeon Phi cores)

Installed August 2013

20 departments

51 faculty (as of Aug. 2014)

#39 on June 2014 Top 500

- Research computing has historically provided some storage for research data for HPC users:

  - Archive (Fortress)

  - Actively running jobs (Cluster Scratch - Lustre)

  - Home directories

**… And Purdue researchers have PURR to package, publish, and describe research data.**

PURDUE
UNIVERSITY

## Scratch needs are climbing

**Avg GB used (top 50 users)**



A bar chart titled "Avg GB used (top 50 users)" with y-axis from 0 to 8000 and the following categories: Steele (2008), Coates (2009), Rossmann (2010), Hansen (2011), Carter (2012), Conte (2013). Legend: Avg GB used (top 50 users).

PURDUE UNIVERSITY

## Fortress usage is skyrocketing

### Fortress Archive Growth



- 2,573,068.00
- 1,002,297.00
- 550,146.00
- 110,909.00

Series1

PURDUE
UNIVERSITY

# HPC STORAGE

Scratch Storage

Fast, large, purged, coupled with clusters, per-user – **for running jobs**

Working Space

Medium speed, large, persistent, data protected, purchased, per research lab – **for shared data and apps**

Archival Space

High-speed, infinite capacity, highly-protected, available to all researchers – **for permanent storage**

PURDUE
UNIVERSITY

# HPC STORAGE

| | $HOME | /group/... | $RCAC_SCRATCH | /tmp | Fortress (HPSS) |
|---|---|---|---|---|---|
| **Capacity** | 10-100 GB | 500 GB and up | Varies by cluster... 500 GB and up | 150-400 GB | unlimited |
| **Resilience to hardware failures** | yes | yes | yes | no | yes |
| **Resilience to human errors** | yes (snapshots) | yes (snapshots) | no | no | no |
| **Subject to purging** | no | no | yes | yes | no |
| **Performance** | medium | medium | high | medium to slow (Hansen) | slow to very slow |
| **Designed for HPC (running jobs off it)** | no | no | yes | no | -No (as main I/O) -Yes (for staging and archiving) |
| **Common access within cluster** | yes | yes | yes | no | yes (hsi/htar) |
| **Common access across clusters** | yes | yes | no (except front-ends) | no | yes (hsi/htar) |
| **Advanced ACLs (beyond ugo/rwx)** | no | yes | no | no | no |

PURDUE
UNIVERSITY

**Working with other researchers across campus, we encounter many different data solutions..**

**From something at the department/workgroup level:**



PURDUE
UNIVERSITY

# RESEARCH DATA ACROSS CAMPUS

## To This





PURDUE
UNIVERSITY

**And This**

## PI Interview Responses:
## How do You Handle Data Storage/Backup?

**Common Data Storage Devices and Services Utilized**

Server hard drives, 65%

Other, 40%

Amazon S3, 8%

Dropbox, 33%

SDSC Project Storage, 13%

SDSC Cloud, 8%

Tape Library, 5%

SAN Storage Array, 3%

USB Drive, 70%

Network Accessible Storage, 70%

*Numbers reflect percentages of PIs surveyed that utilize each solution ; Individual PIs use multiple solutions, so %'s add up to >100%.*

- **Storage Devices**
  - Network accessible storage (NAS), USB and server local drives dominate
  - Use of Dropbox for sharing
  - Others use Google Drive, Hadoop, XSEDE, SDSC co-location
- **Backup modes**
  - Replicated copies in two NAS
  - A copy in the NAS,
  - A copy in local hard drive (laptop/workstation),
  - And a copy in a USB drive
  - Maybe a copy in email/Dropbox
- **Problems:**
  - Out of sync
  - Lost track of its location
  - Lost version control
  - High cost of recovery

**UNIVERSITY OF CALIFORNIA, SAN DIEGO**

PURDUE UNIVERSITY

RCi UC San Diego Research Cyberinfrastructure

UCSD

# RESEARCH STORAGE GAPS

**Many central storage options have not met all the needs that researchers care about**

- Departmental or lab resources are islands and not accessible from HPC clusters.

- Most are individually-oriented, rather than built around the notion of a research lab.

  - Boiler Backpack, Fortress, research homes, and scratch

- Scratch filesystems are *also* limited in scope to a single HPC system

PURDUE
UNIVERSITY

# RESEARCH STORAGE GAPS

**Before 2013, we've heard lots of common requests:**

- I need more space than I can get in scratch

- Where can I install applications for my entire research lab?

- I'm *actively working* on that data/software in scratch:

    - I have to go to great lengths to keep it from being purged.

    - I shouldn't have to pull I from Fortress over and over

- Can I get a UNIX group created for my students and I?

- Is there storage that I can get to on *all* the clusters I use?

- I have funding to spend on storage – what do you have to sell?

- I need storage for my instrument to write data into

- My student has the only copy of my research data in his home directory, and he graduated/went off the grid!

PURDUE
UNIVERSITY

**We've addressed some of these with improving scratch:**

- Everybody automatically gets access to Fortress for permanent data storage

- Very large per-user quotas, beginning on Conte

- More friendly purge policy – based on the *use* of data, rather than creation time.

# PERSISTENT

## GROUP STORAGE

### A STORAGE SERVICE FOR HPC USERS

# A SOLUTION

**Since early 2013, HPC researchers could at last purchase storage!**

Quickly, a group storage service for research began to materialize to address many common requests:

- 500G available at no additional charge to community cluster groups

- Mounted on all clusters and exported via CIFS to labs

- *Not scratch*: Backed up via snapshots

- Data in /group is owned by faculty member!

- Sharing ability – Globus, CIFS, and WWW

- Version Control repositories

- Maintain group-wide copies of application software or shared data

PURDUE UNIVERSITY

```
/group/mylab/

              +--/apps/
              |
              +--/data/
              |
              +--/etc/
              |           +--bashrc
              |           +--cshrc
              +--/repo/
              |
              +--/www/
              |
              +--/...
```

…with POSIX ACLs to overcome group permission and umask challenges!

PURDUE
UNIVERSITY

SELF-SERVICE

MANAGE YOUR OWN ACCESS

# Fairly well received!

- In just over one year, over 65 research groups are participating.

  - *Several are not HPC users!*

- Over .5 PB provisioned to date

- A research group purchasing space has purchased, on average, 8.6TB.

## It's not perfect

- *Scalability and I/O latency are issues on clusters*

- Research computing staff have to monitor "appropriate" work

    - An I/O intensive job in the wrong directory hurts the entire campus!

- Creation and management is not automated

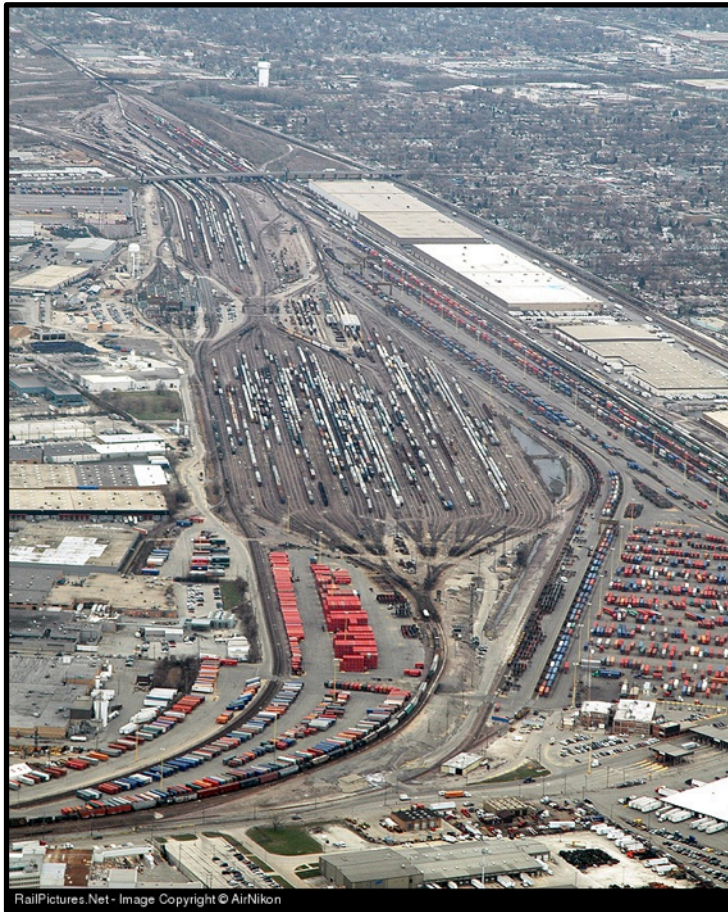- Protected from "oops" – but **not** disaster-protected

# THE RESEARCH
# DATA DEPOT
## INFRASTRUCTURE FOR RESEARCH DATA

PURDUE UNIVERSITY

# DEPOT

As a transport hub: a place where large amounts of cargo are stored, loaded, unloaded, and moved from place to place.



PURDUE
UNIVERSITY

A heavily protected and impenetrable building *(Oxford)*

# DESIGN GOALS

**HPC Faculty were surveyed, we learned that they require from their storage:**

- Protection from accidents (snapshots)

- Protection from disaster (replicas)

- Affordable prices

**We also know that we need to provide:**

- A high-performance resource for research instruments

- A business model to sell storage, if faculty wish to purchase it

- Easy ways to share data with collaborators

- Easier ways to access Fortress

- A resource that works for both HPC *and* non-HPC users

PURDUE
UNIVERSITY

# RFP REQUIREMENTS

## Proposals Requested from Vendors that can deliver:

| Depot Requirements | What we have today |
|---|---|
| At least 1 PB usable capacity | >1 PB |
| 40 GB/sec throughput | **5 GB/sec** |
| < 3ms average latency, < 20 ms maximum latency | **Variable** |
| 100k IOPS sustained | **55k** |
| 300 MB/sec min client speed | **200 MB/sec max** |
| Support 3000 simultaneous clients | Yes |
| Filesystem snapshots | Yes |
| Multi-site replication | **No** |
| Expandable to 10 PB | Yes |
| Fully POSIX compliant, including parallel I/O | **No** |

**PURDUE**
UNIVERSITY

# SOLUTION

Approximately 2.25 PB of IBM GPFS

Hardware provided by a pair of Data Direct Networks SFA12k arrays, one in each of MATH and FREH datacenters

160 Gb/sec to each datacenter

5x Dell R620 servers in each datacenter



PURDUE
UNIVERSITY

**At $150/TB per year:**

- Same base capabilities as "group storage", **PLUS**

  - Snapshots

  - Multi-site copies of your data

  - A scalable, expandable storage resource optimized for HPC

- Access to Globus data transfer service, and endpoint sharing

PURDUE
UNIVERSITY

# RELATED
# CAPABILITIES
## EVERYTHING IS ABOUT DATA!

## 2014 network improvements

- 100 Gb/sec WAN connections

- Research Core

    - 160 Gb/sec core to each resource (up from 40)

    - 20 Gb/sec research core to most of campus

- Campus Core Upgrade

https://www.rcac.purdue.edu/news/681



Internet2 Combined Infrastructure Topology

PURDUE UNIVERSITY

## Globus

- Move data between Purdue, other institutions, and national HPC resources

- Easily share data with collaborators around the world
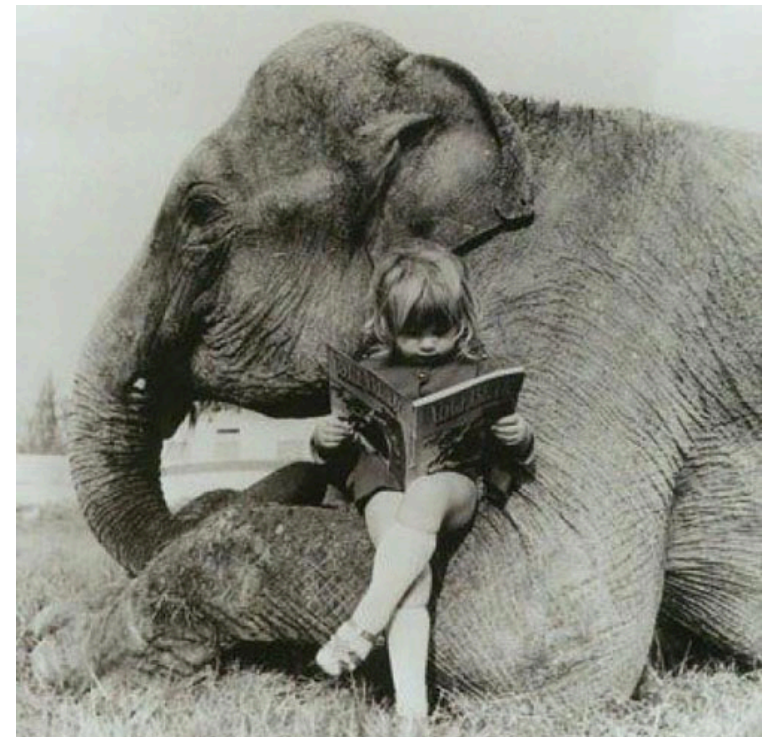
  - With dropbox-like simplicity!



*https://transfer.rcac.purdue.edu*

PURDUE
UNIVERSITY

# BIG DATA

- Use the new "hathi" Hadoop cluster for prototyping big data applications

- Spark, Hbase, Hive, Pig

- On-demand *MyHadoop* clusters on your group's Conte or Carter nodes

**https://www.rcac.purdue.edu/compute/hathi/**

more awesome pictures at THEMETAPICTURE.COM

## Firebox virtual servers

- Host LAMP servers, cluster login nodes, submission portals, nontraditional HPC, or interactive desktops, all within the research infrastructure

- Move your computing environment close to your research data!



https://www.rcac.purdue.edu/services/firebox/



**PURDUE**
UNIVERSITY

**Group storage in Fortress**

- Your same self-managed group for the Depot can manage access to a group fortress space

- Contact RCAC for assistance with setting this up



PURDUE UNIVERSITY

# THE RESEARCH

## DATA DEPOT

### WHAT IS NEXT?

**The Research Data Depot is NOT:**

- A 1:1 replacement for Dropbox

- Blessed for HIPAA, FISMA, ITAR, or other regulated data

- Global, paid-for, or non-purged scratch

  - Each resource will continue to have its own local, dedicated, high-performance scratch

  - Its *is* designed to not **actively hurt** an HPC user or others should it be used for running jobs

- A data publishing service

PURDUE
UNIVERSITY

**Potential future questions**

- How can we more tightly integrate the Depot and Fortress?

- How can we bridge the space between the Depot at PURR?

- Can your Depot-using labs and instruments benefit from the campus' improved networking infrastructure?

PURDUE
UNIVERSITY

**Recent or in-development information sessions on working with data**

- Sept 2: "Big Data" training session

- Globus

- Effective use of HPC storage

- HPC I/O formats: HDF5

- In collaboration with Purdue Libraries – the data life cycle and managing your research data on RCAC resources

PURDUE
UNIVERSITY

- Collaborations on multi-disciplinary grant proposals, both internal and external

- Developing customized Data Management Plans

- Organizing your data

- Describing your data

- Sharing your data

- Publishing your datasets

- Preserving your data

- Education on data management best practices

**PURDUE**
UNIVERSITY

# LIBRARY DATA SERVICES

**PURR** is a *free* online research data collaboration platform and service solution for Purdue faculty, graduates students, and staff.

**Research data** - spreadsheets, images, output from sensors and instruments, transcripts, surveys, software source code and tools, video, and observation logs

**PURR provides:**

✓ Data management plan (DMP) resources and consultation

✓ Collaborative research data project space

✓ Dataset publication with Digital Object Identifier (DOI) *

✓ Long-term preservation and management

PURDUE
U N I V E R S I T Y

## HOW CAN I GET ACCESS?

The service is in early access mode.

Contact us at rcac-help@purdue.edu if you're willing to be an early access tester!

**Full production on October 31**

PURDUE
UNIVERSITY

# /GROUP USERS

Persistent group storage owners will be transitioned into the Depot at an agreed-upon date, on a group-by-group basis.

We'll be in touch to arrange a time!

**If you're paid up with the persistent group storage, this move will come at no additional charge until your renewal is due.**

PURDUE
UNIVERSITY

# We can set your lab up with a small evaluation space upon request

## Need 1TB or more?

purchase access online

- http://www.rcac.purdue.edu/order



**PURDUE UNIVERSITY**

# CLUSTER COMPUTING

## WHAT IS NEXT?

# TODAY

**Carter & Conte still very current**

- Ivy Bridge chips not much improvement over Carter
- Conte's Xeon Phi co-processors extremely capable for certain applications (dense FP simulation)
- Both Carter and Conte have proven very popular!



PURDUE
UNIVERSITY

**Released last week:**

- More cores/chip
- Slight clock speed reduction
- Additional floating-point instructions
  - Can greatly assist matrix multiplication operations
- DDR4 memory
  - Higher memory bandwidth helps just about everyone

**Longer-term implications:**

- Forget 2/4/8 GB/core – think about base memory/node
- Best programming model may be MPI+OpenMPI
  - Fewer MPI ranks/node, pick up extra cores/threads with OpenMP
  - Looks vaguely like Xeon Phi…

**PURDUE**
U N I V E R S I T Y

**Planning for next community cluster starting now**

- Bringing Haswell-based loaner to Purdue

- Preparing for open bid later this fall

- Cluster assembly, go-live in mid-Spring, 2015

**Open questions:**

- Higher clock speed or more cores/node?

- Memory sizes?
    - Probably 64 GB base with an option for large (256G or more) nodes

PURDUE
U N I V E R S I T Y

# THE END

Questions?

PURDUE
UNIVERSITY