# Accessing Multi-TB-sized Datasets at NCAR's Research Data Archive

## Douglas Schuster, Steven Worley
Database Engineer, CISL/DSS
http://rda.ucar.edu

## ECMWF 16th MOS Workshop,  Reading, UK
March 1, 2017

# Highlights

- Overview of NCAR and the RDA
- Data Management Strategies
- Metadata Overview
- RDA User Services
- Usage Metrics
- User Outreach and User Support
- Lessons Learned and Future Directions

# Highlights

- **Overview of NCAR and the RDA**
- Data Management Strategies
- Metadata Overview
- RDA User Services
- Usage Metrics
- User Outreach and User Support
- Lessons Learned and Future Directions

# Overview of NCAR

**National Center for Atmospheric Research -1960**

**Mission Statement:**
- To understand the behavior of the atmosphere and related Earth and geospace systems

- To support, enhance, and extend the capabilities of the university community and the broader scientific community, nationally and internationally, and

- To foster the transfer of knowledge and technology for the betterment of life on Earth

# Overview of NCAR

- Federally Funded Research and Development Center
  - National Science Foundation
  - Other government agencies, other national governments, private sector
- Managed by University Corporation for Atmospheric Research (UCAR)
- 5 campuses, 1200+ employees.

# Overview of the RDA



- **NCAR's Research Data Archive**
- **Purposes**
  - Support climate & weather research at NCAR and UCAR universities with reference datasets
- **Collections**
  - Ocean & atmospheric observations, analyses, reanalyses, operational NWP products
  - Established in 1960s
  - 700+ datasets, 8M files, 1.8 PB
  - 70+ datasets growing daily-monthly
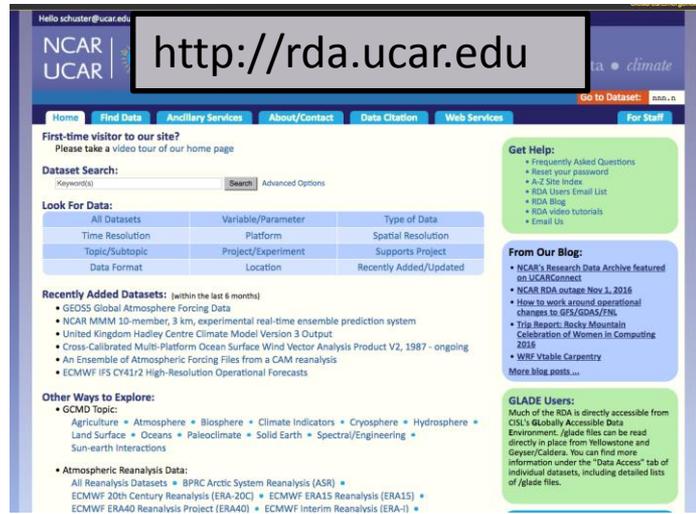- **Free and open access**
- **Worldwide usage**
- **Science educated staff**
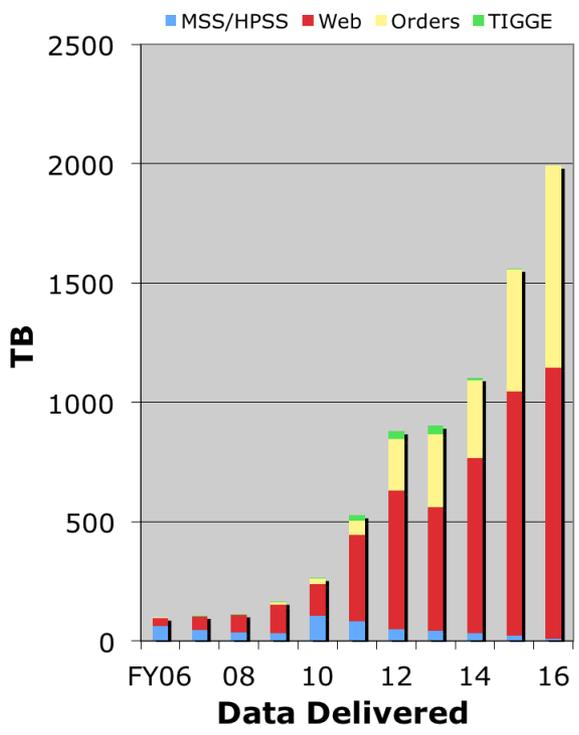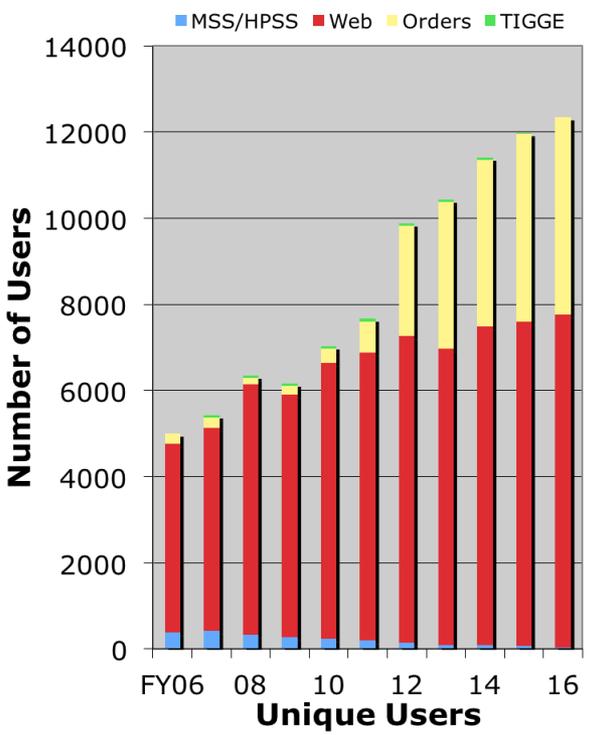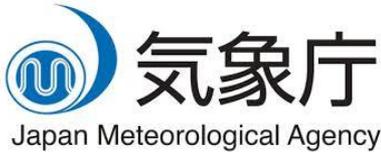- **One of many data assets at NCAR**

# Overview of the RDA



12,000+ Unique Web users in FY 2016

~2 PB Data Delivered in FY 2016

http://rda.ucar.edu

# Overview of the RDA Collaboration Examples



- Host Datasets and Share Data
  - Operational model output
  - Reanalysis products
  - Observations to support reanalysis efforts
- Develop Products and Services
  - E.G. ICOADS ,ISPD, and NCAR UA soundings
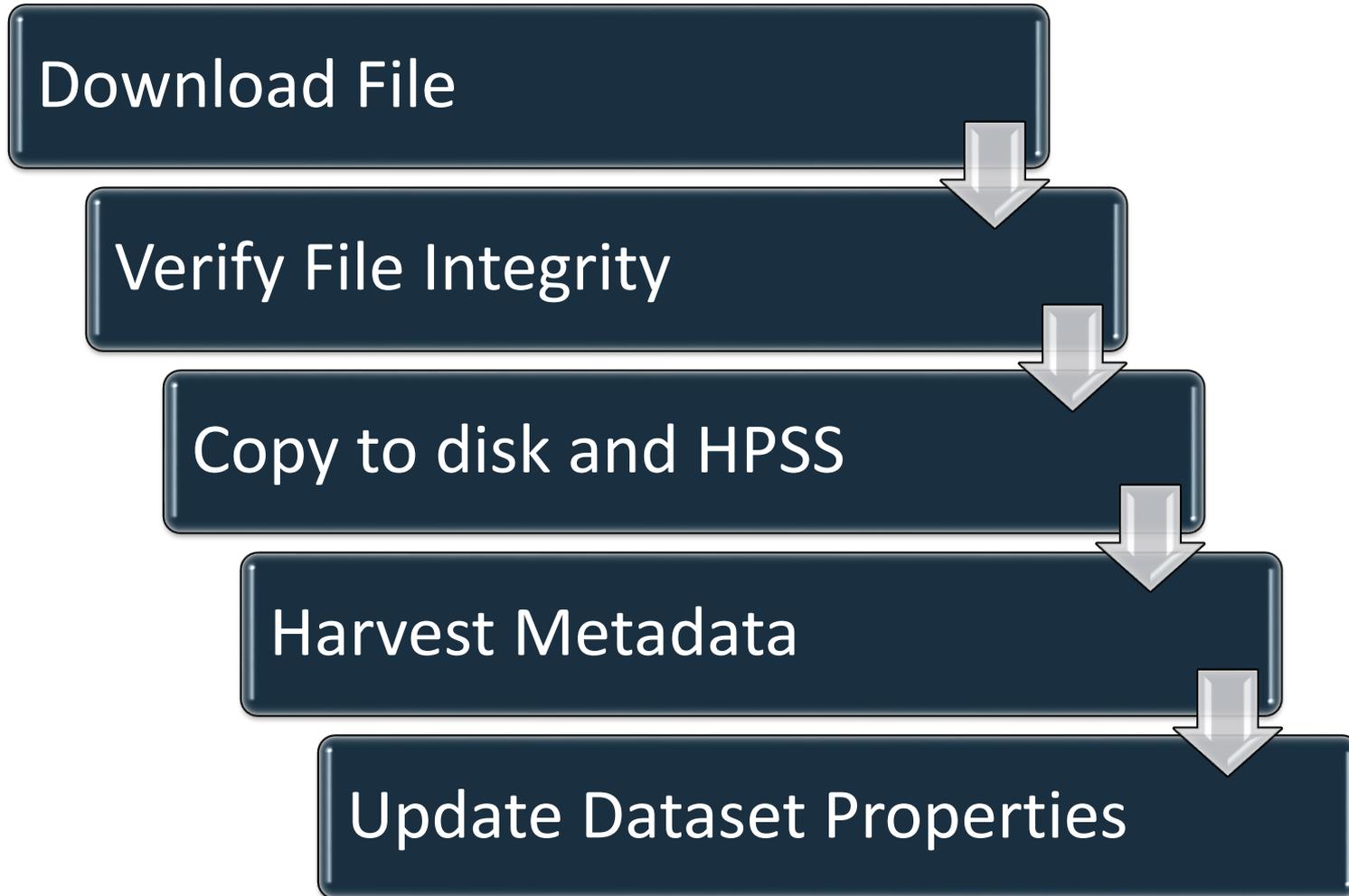  - Interpolate and restructure products

# Highlights

- Overview of NCAR and the RDA
- Data Management Strategies
- Metadata Overview
- RDA User Services
- Usage Metrics
- User Outreach and User Support
- Lessons Learned and Future Directions

# Data Management Strategies

- Develop and leverage tools that support programmatic data ingest
  - Support scalable archive growth and maintenance
- Establish and maintain metadata databases
  - Support discovery, access, and value added services
- Store data at multiple physical locations
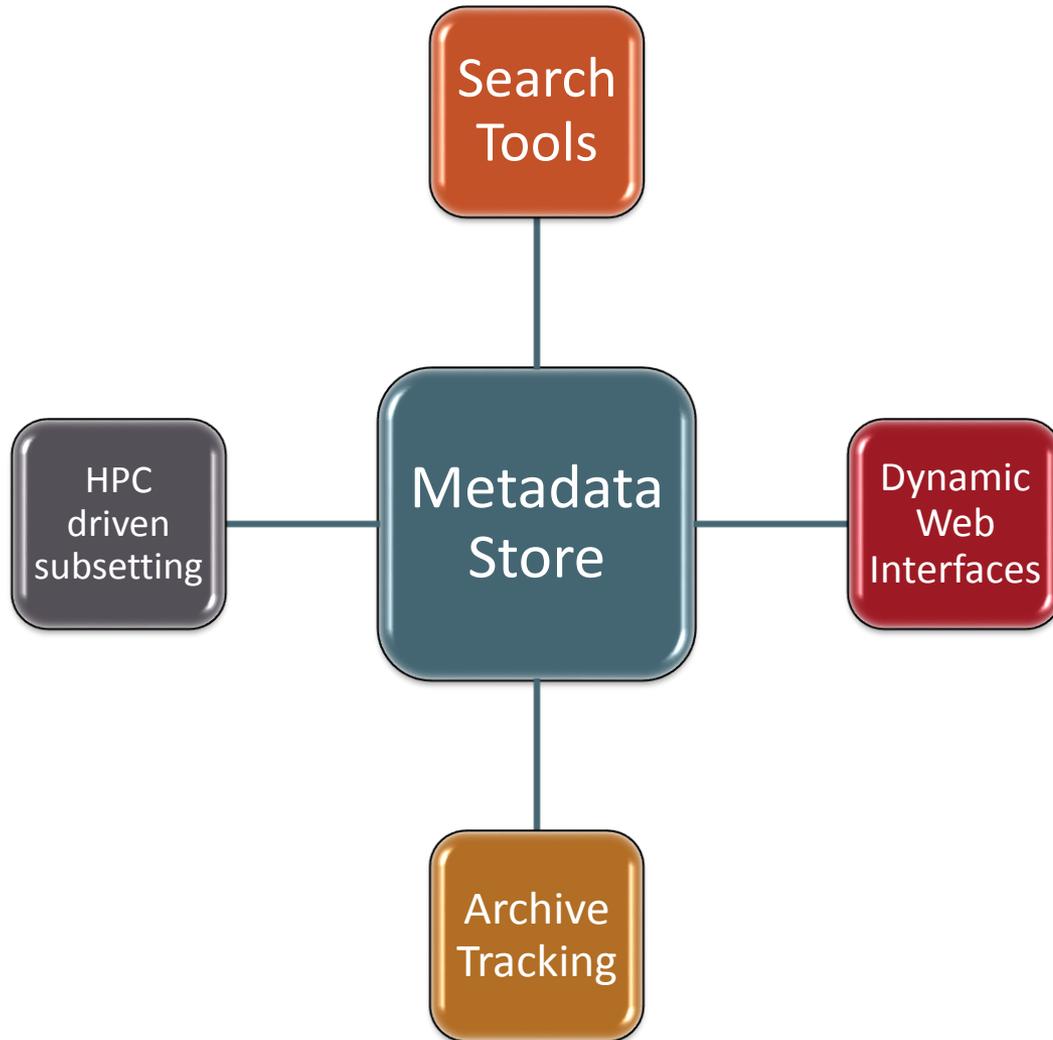  - Support data security

# Data Management Strategies
# Typical Dataset Update Workflow

Download File

Verify File Integrity

Copy to disk and HPSS

Harvest Metadata

Update Dataset Properties

# Data Management Strategies
# Metadata is foundational

# Highlights

- Overview of NCAR and the RDA
- Data Management Strategies
- **Metadata Overview**
- RDA Users Services
- Usage Metrics
- User Outreach and User Support
- Lessons Learned and Future Directions

# Metadata Overview

- **Metadata**
  - Dataset summary metadata stored in native RDA schema with GCMD control vocabulary
  - Exported to
    - ISO 19139, DIF, FGDC, OAI_DC, DATACITE, THREDDS
  - Shared via
    - OAI-PMH endpoint
    - CS/W endpoint
  - Harvested by
    - NCAR data search and discovery system
    - NASA Global Change Master Directory (GCMD)
    - Thompson Reuters
    - GEOSS –GEO portal

# Highlights

- Overview of NCAR and the RDA
- Data Management Strategies
- Metadata Overview
- RDA User Services
- Usage Metrics
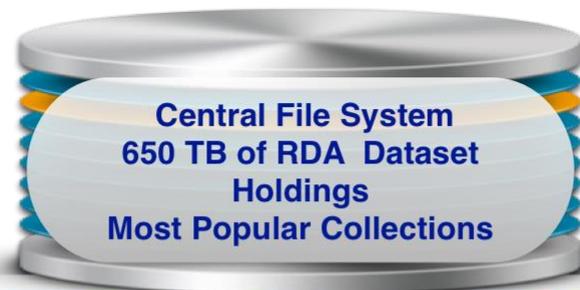- User Outreach and User Support
- Lessons Learned and Future Directions

# RDA User Services

- Make data discovery "easy"
- Improve access to serve growth in user numbers and types
- Support reproducible research
  - DOI's on 79 static and dynamic datasets
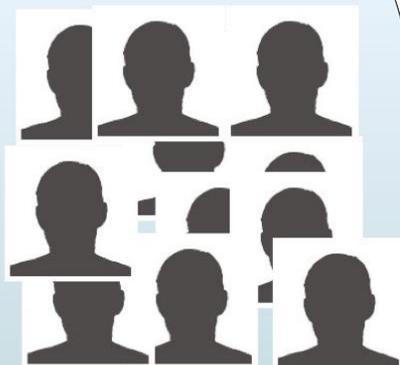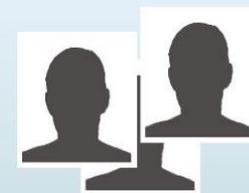- **Reduce the time researchers spend dealing with data**

# RDA User Services



HPSS Tape Library
1.2 PB of RDA Dataset Holdings
Full Copy of Archive

# RDA User Services

**Central File System**
**650 TB of RDA Dataset Holdings**
**Most Popular Collections**

**HPSS Tape Library**
**1.2 PB of RDA Dataset Holdings**
**Full Copy of Archive**

# RDA User Services



**External Users ~12K/year**

**NCAR HPC Users ~100s/year**

**Central File System**
**650 TB of RDA Dataset Holdings**
**Most Popular Collections**

**HPSS Tape Library**
**1.2 PB of RDA Dataset Holdings**
**Full Copy of Archive**

# RDA User Services



External Users ~12K/year
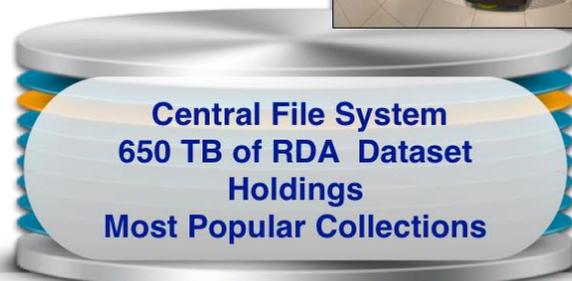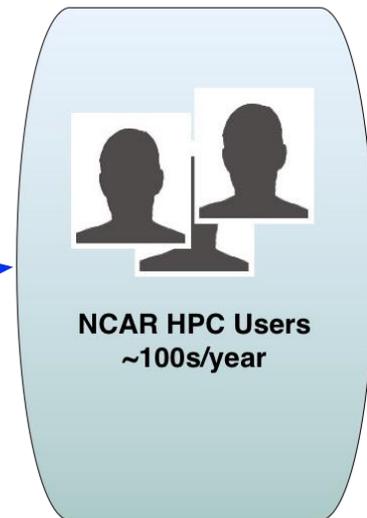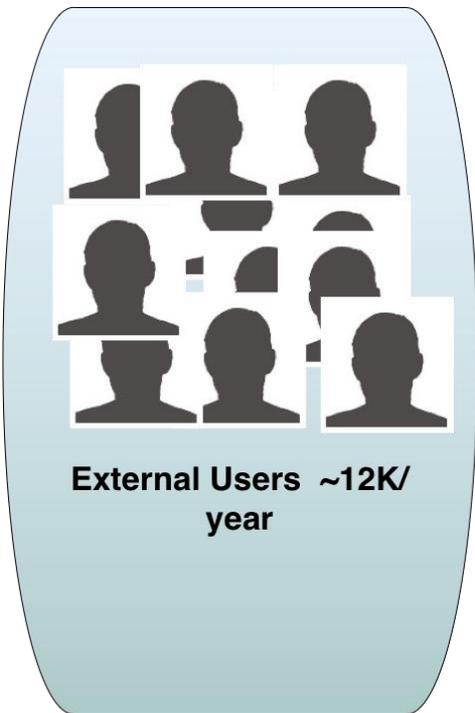
High Performance Computing

NCAR HPC Users ~100s/year

Central File System
650 TB of RDA Dataset Holdings
Most Popular Collections

HPSS Tape Library
1.2 PB of RDA Dataset Holdings
Full Copy of Archive

# RDA User Services



**External Users ~12K/year**

**Web Servers**

**High Performance Computing**

**NCAR HPC Users ~100s/year**

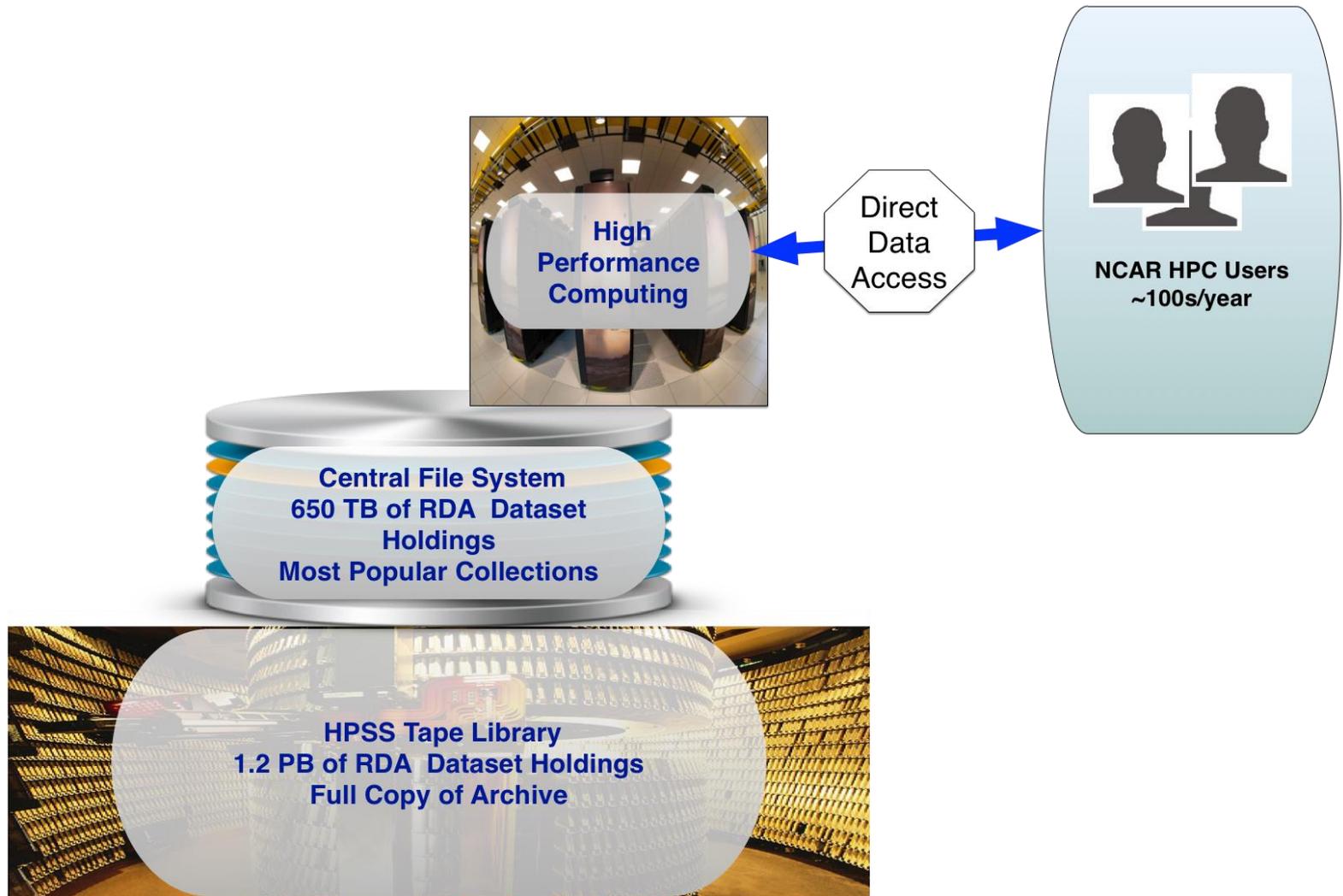**Central File System**
**650 TB of RDA Dataset Holdings**
**Most Popular Collections**

**HPSS Tape Library**
**1.2 PB of RDA Dataset Holdings**
**Full Copy of Archive**

# RDA User Services



**High Performance Computing**

Direct Data Access

**NCAR HPC Users ~100s/year**

**Central File System**
**650 TB of RDA Dataset Holdings**
**Most Popular Collections**

**HPSS Tape Library**
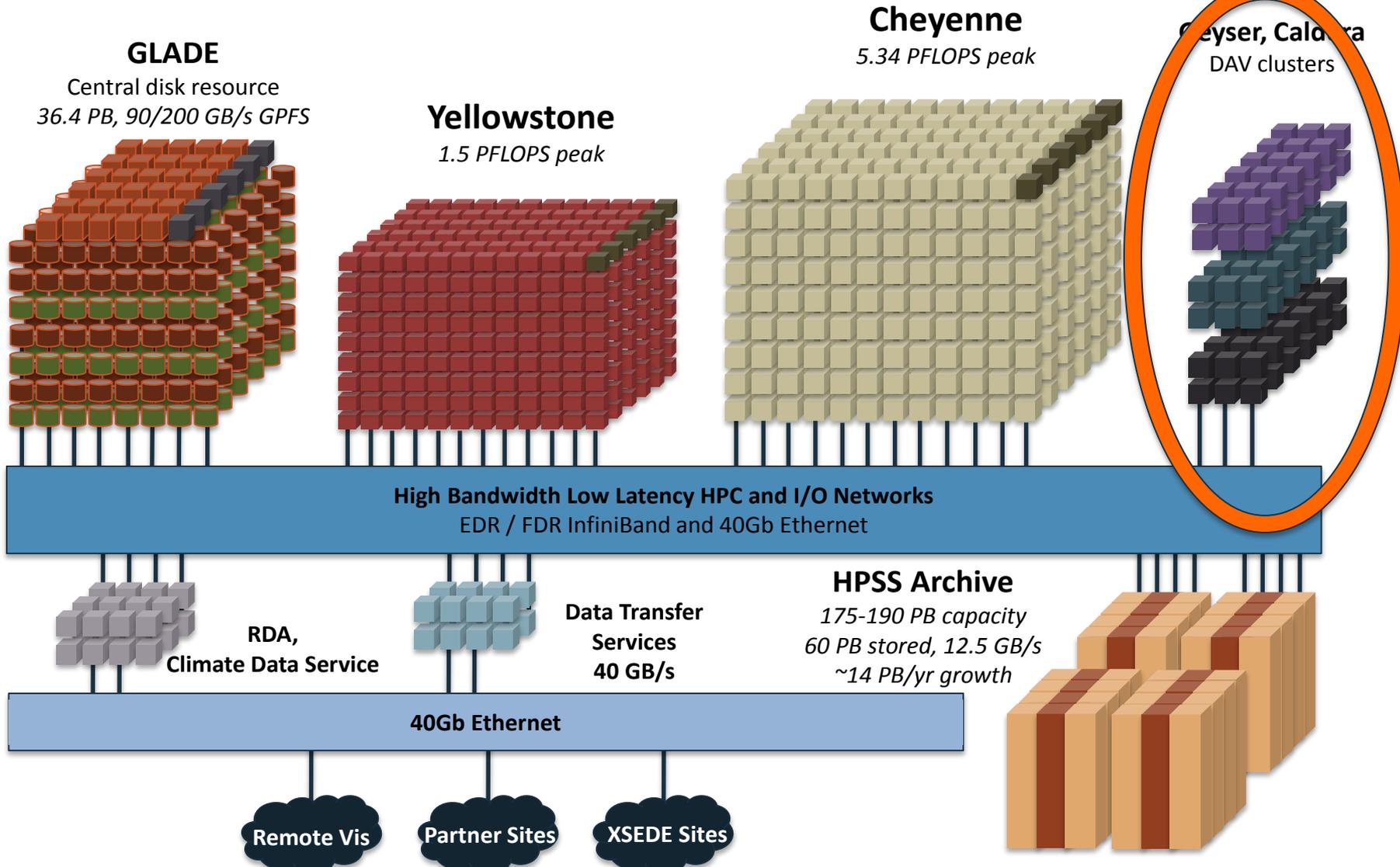**1.2 PB of RDA Dataset Holdings**
**Full Copy of Archive**
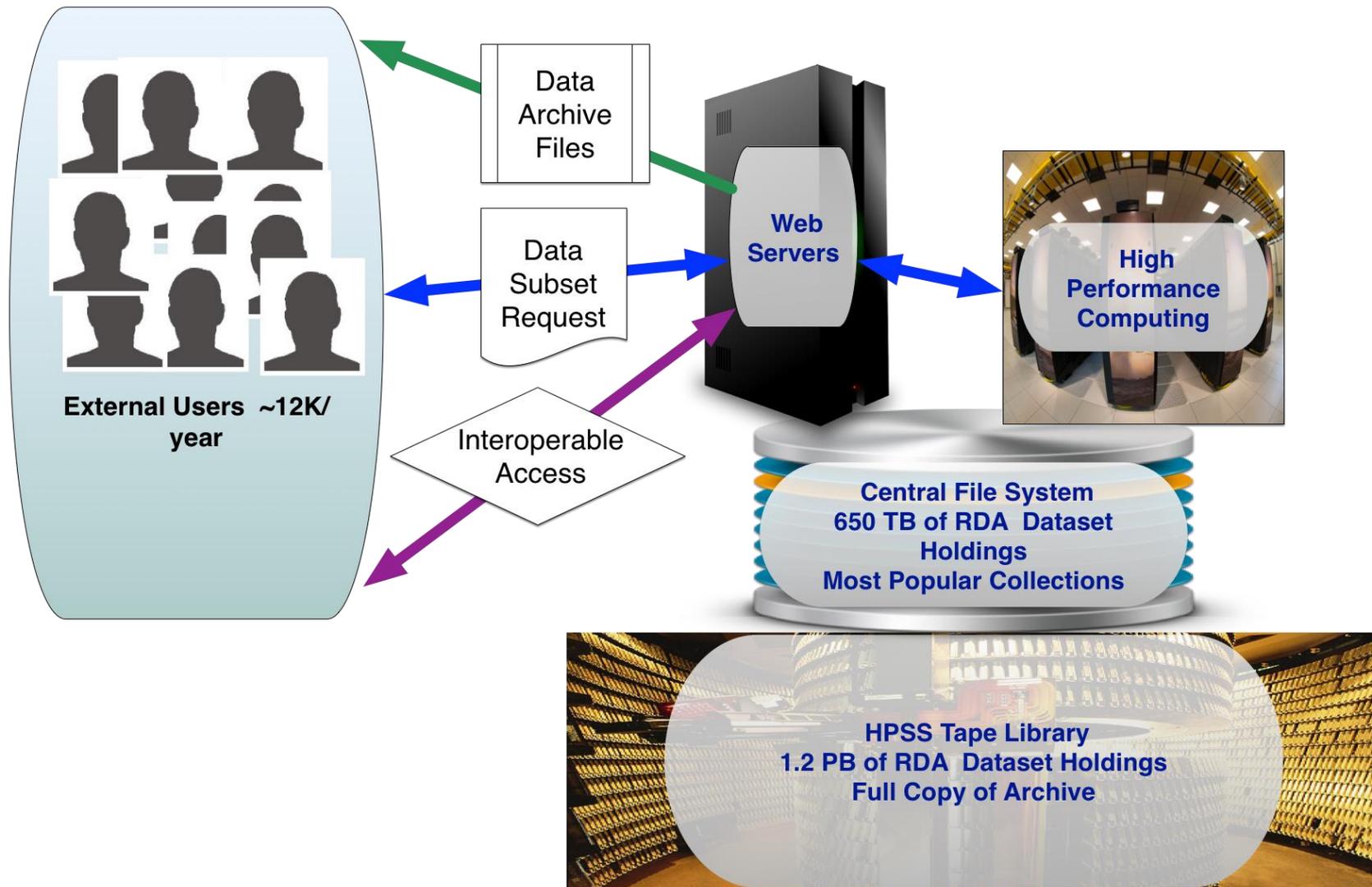
# NCAR Wyoming Supercomputing Center

- Opened in 2012
- Building performing great
- Ongoing improvements to achieve best operational efficiency
- Cheyenne (SGI) system brought online in Feb 2017 -#20 in top 500
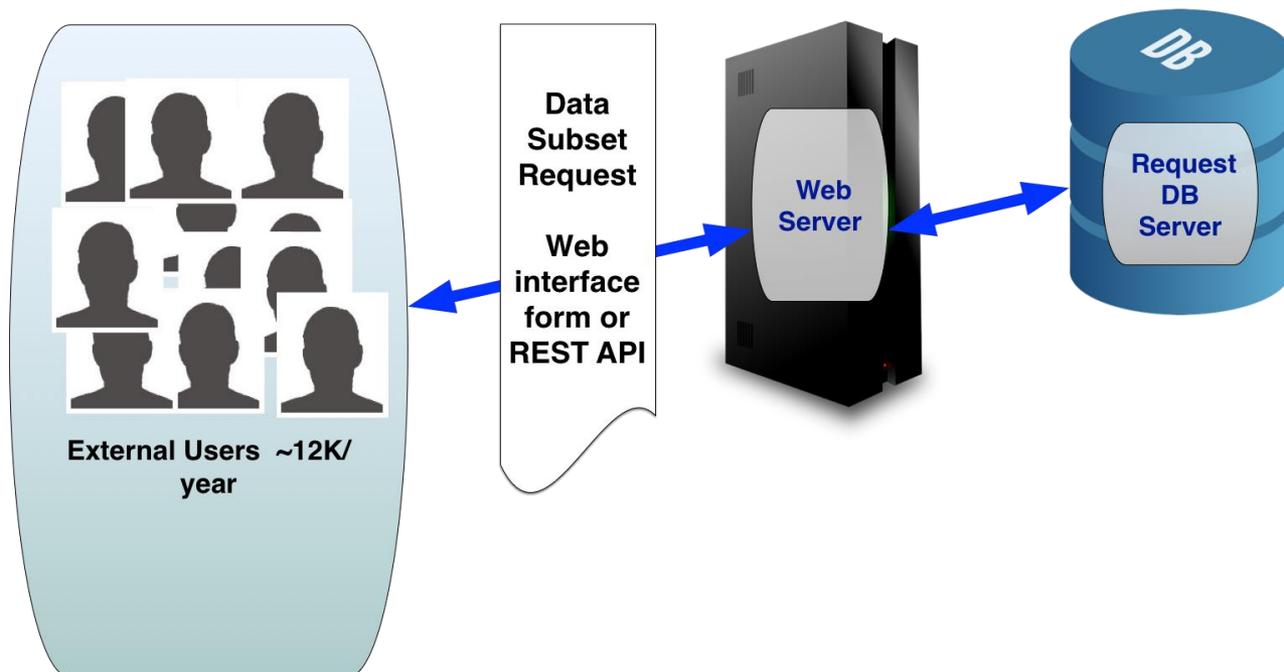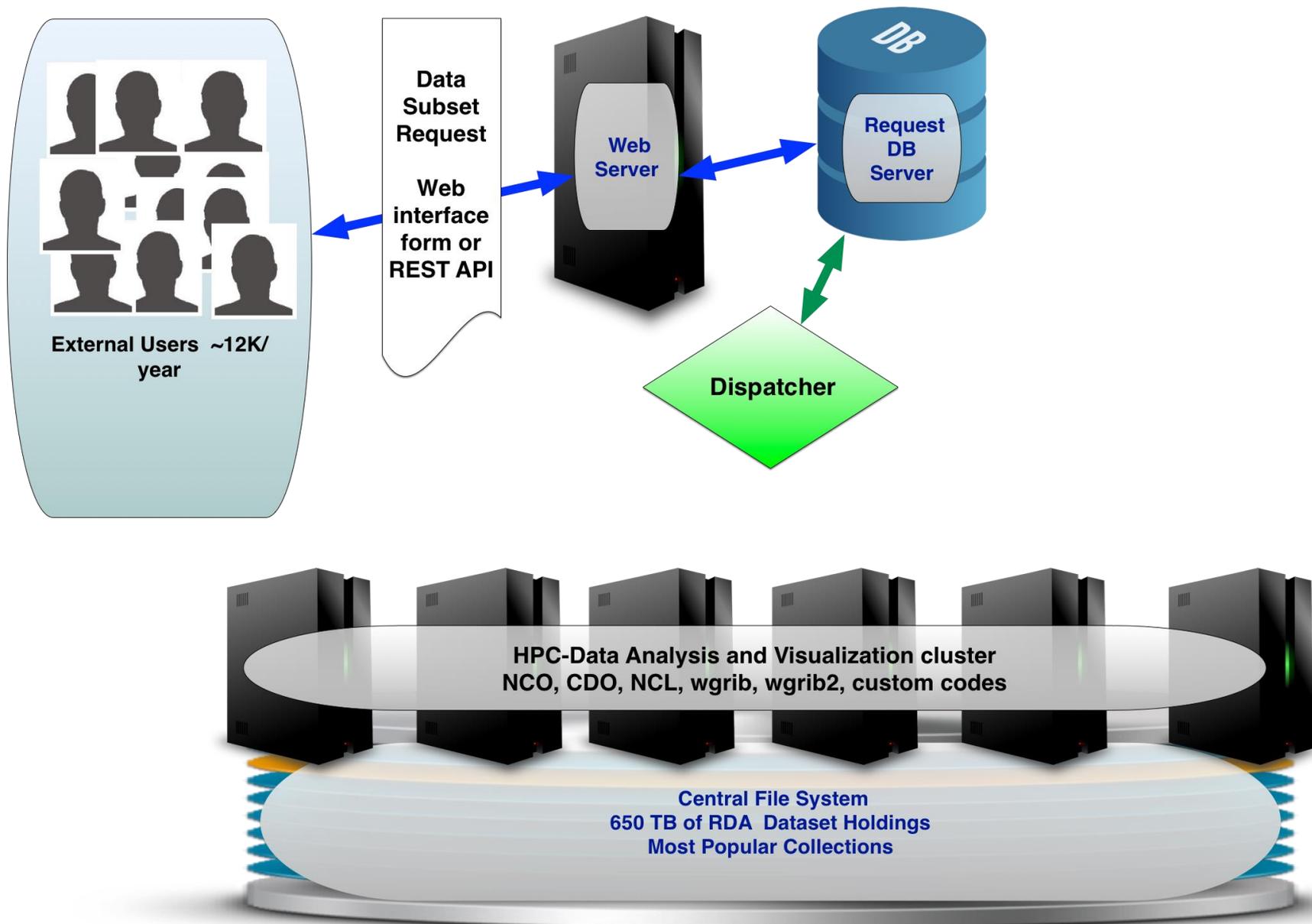
# NCAR Data Intensive Supercomputing Environment

**GLADE**
Central disk resource
*36.4 PB, 90/200 GB/s GPFS*

**Cheyenne**
*5.34 PFLOPS peak*

**Geyser, Caldera**
DAV clusters

**Yellowstone**
*1.5 PFLOPS peak*

**High Bandwidth Low Latency HPC and I/O Networks**
EDR / FDR InfiniBand and 40Gb Ethernet

**RDA,
Climate Data Service**

**Data Transfer
Services
40 GB/s**

**HPSS Archive**
*175-190 PB capacity
60 PB stored, 12.5 GB/s
~14 PB/yr growth*

**40Gb Ethernet**

**Remote Vis**

**Partner Sites**

**XSEDE Sites**

# RDA User Services



Data Archive Files

Data Subset Request

Interoperable Access

Web Servers

High Performance Computing

External Users ~12K/year

Central File System
650 TB of RDA Dataset Holdings
Most Popular Collections

HPSS Tape Library
1.2 PB of RDA Dataset Holdings
Full Copy of Archive

# RDA User Services
# Data Volume Reduction through HPC Processing



External Users ~12K/year

Data Subset Request

Web interface form or REST API

Web Server

Request DB Server

HPC-Data Analysis and Visualization cluster
NCO, CDO, NCL, wgrib, wgrib2, custom codes

Central File System
650 TB of RDA Dataset Holdings
Most Popular Collections

# RDA User Services
# Data Volume Reduction through HPC Processing

Data
Subset
Request

Web
interface
form or
REST API

Web
Server

Request
DB
Server

DB

Dispatcher

External Users ~12K/ year

HPC-Data Analysis and Visualization cluster
NCO, CDO, NCL, wgrib, wgrib2, custom codes

Central File System
650 TB of RDA Dataset Holdings
Most Popular Collections

# RDA User Services
# Data Volume Reduction through HPC Processing



External Users ~12K/year

Data Subset Request

Web interface form or REST API

Web Server

Request DB Server

DB

Dispatcher

Files 1-100

Files 101-200

Files 201-300

Files 301-400

Files 401-500

Files 501-600

**HPC-Data Analysis and Visualization cluster
NCO, CDO, NCL, wgrib, wgrib2, custom codes**

**Central File System
650 TB of RDA Dataset Holdings
Most Popular Collections**

# RDA User Services
# Data Transfer Mechanisms



**External Users ~12K/ year**

http://

globus

**GridFTP Servers**

**Web Server**

**Central File System**
**150TB of cache to hold requests -5 day retention default**

# RDA User Services
# Data Transfer Mechanisms
# Globus Advantages



- **Reliable, secure, high-performance** file transfer
- "Fire and forget"
- Automatic fault recovery
- Powerful GUI, APIs, and CLI



transfer

secure endpoint

**2** Globus moves the data for you

secure endpoint

A

B

**1** You submit a transfer request

**3** Globus notifies you once the transfer is complete

# RDA User Services
# OpenDAP Access Example Product, THREDDS



Specific humidity @ Isobaric Surface

GRIB reference time: 1975-06-01 00:00:00 +0000

Specific humidity @ Isobaric Surface (kg kg^-1)

1.2E-05    3.2E-03    6.3E-03    9.4E-03    1.3E-02    1.6E-02

Data Min = 1.2E-05, Max = 1.6E-02, Mean = 6.3E-03

# Highlights

- Overview of NCAR and the RDA
- Data Management Strategies
- Metadata Overview
- RDA User Services
- **Usage Metrics**
- User Outreach and User Support
- Lessons Learned and Future Directions

# Current Services Overview – Usage metrics



Traditional File Download

User specified extraction

# Traditional File Download – Usage metrics



**Yearly RDA User Access from Web Interface**
Direct Archive File Downloads

In FY 2016
- 7900 users
- 1.2 PB of data delivered

# Breaking the Big Data barrier for users

- **Users select what they want, RDA extracts what they need**

## Users, Data Accessed & Exported

|  | 2012 | 2013 | 2014 | 2015 | 2016 |
|---|---|---|---|---|---|
| Users | 2600 | 3400 | 3900 | 4400 | 4600 |



% Exported

% Accessed

Pie area proportional to # requests

| # Requests | 14000 | 25000 | 44000 | 41000 | 76000 |
|---|---|---|---|---|---|
| PB Access | 6 | 14 | 16 | 22 | 30 |
| TB Export | 226 | 315 | 332 | 534 | 846 |
| % Reduced | 96 | 98 | 98 | 97 | 97 |

# Highlights

- Overview of NCAR and the RDA
- Data Management Strategies
- Metadata Overview
- RDA User Services
- Usage Metrics
- User Outreach and User Support
- Lessons Learned and Future Directions

# User Outreach and Support
# -Server Generated Data Citation



**How to Cite This Dataset:**

**RIS**
**BibTeX**

Compo, G. P., and Coauthors, 2009: NOAA CIRES Twentieth Century Global Reanalysis Version 2. Research Data Archive at the National Center for Atmospheric Research, Computational and Information Systems Laboratory, Boulder, CO. [Available online at http://dx.doi.org/10.5065/D6QR4V37.] Accessed† dd mmm yyyy.

†Please fill in the "Accessed" date with the day, month, and year (e.g. - 5 Aug 2011) you last accessed the data from the RDA.

Bibliographic citation shown in  American Meteorological Society (AMS)  style

Get a customized data citation

**Dataset:**

NOAA CIRES Twentieth Century Global Reanalysis Version 2 (ds131.1)

**Your Access History for July 2013:**

Choose a day to get a citation for this dataset. You will also see more details about your downloads on that day, which will help you verify that this is the citation you want.

| July 2013 | | | | | | |
|---|---|---|---|---|---|---|
| Sun | Mon | Tue | Wed | Thu | Fri | Sat |
|  | 1 | 2 | **3** | 4 | 5 | 6 |
| 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| 14 | 15 | 16 | 17 | 18 | 19 | 20 |
| 21 | 22 | 23 | 24 | 25 | 26 | 27 |
| 28 | 29 | 30 | 31 |  |  |  |

Maintain user access history, citation recall at any time

**For Data Accessed on 2013-07-03:**

**Dataset Citation:**   **RIS**

Compo, G. P., et al. 2009. *NOAA CIRES Twentieth Century Global Reanalysis Version 2*. Research Data Archive at the National Center for Atmospheric Research, Computational and Information Systems Laboratory. http://dx.doi.org/10.5065/D6QR4V37. Accessed 3 Jul 2013.

Bibliographic citation shown in  Federation of Earth Science Information Partners (ESIP)  style

**Data Access Detail:**

1 subset request:

- 18 files, 3.44 MB
  - Date Limits      :  1975-06-15 00:00 to 1975-06-30 12:00
  - Parameter        :  TMP
  - Level Type       :  HGT

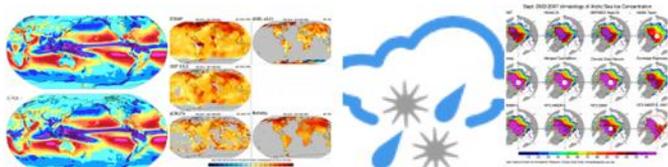# User Outreach and Support -Advise User of Best Dataset Collection(s) for purpose

# User Outreach and Support
# Social Media, with User Video Tutorials

# User Outreach and Support Direct Help

- Dataset consultant listed on each dataset homepage



- General help email rdahelp@ucar.edu

# Highlights

- Overview of NCAR and the RDA
- Data Management Strategies
- Metadata Overview
- RDA User Services
- Usage Metrics
- User Outreach and User Support
- **Lessons Learned and Future Directions**

# Lessons Learned

- Archive structure impacts efficiency of data access services
  - Use cases: climate vs weather research
  - Use of tar packages/compression
- Programmatic metadata harvesting critical to support value added services, and data integrity validation upon ingest.
- Essential to co-locate data with HPC on disk

# Future Directions

- CAPSTONE Data Analysis Environment
  - Develop and share data analysis workflows
- SSD Storage to drive fast access
- Collaborate with NCAR GIS program
- RDA-ODB instance –Support NCAR DA research
- Distributed Ocean Matchup Service (DOMS)
  - Satellite to in-situ matchup distributed archives
- Continue to add high value assets to archive
  - ERA-5, ERA-CLIM
  - ASR
  - NCAR Reanalysis Products

Questions?

http://rda.ucar.edu

Doug Schuster
schuster@ucar.edu