# *Large-scale Data Systems at the National Center for Atmospheric Research (NCAR)*
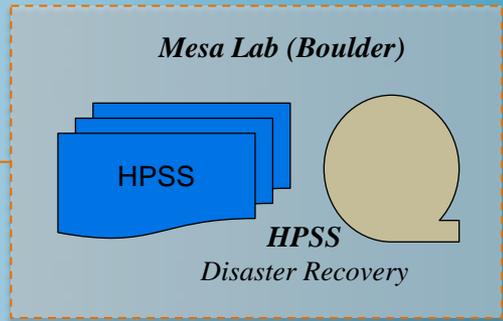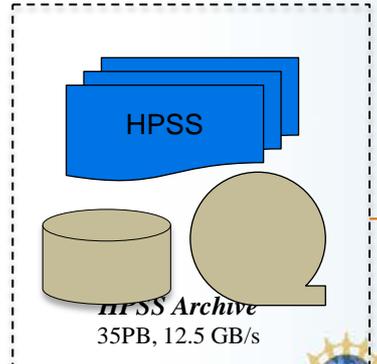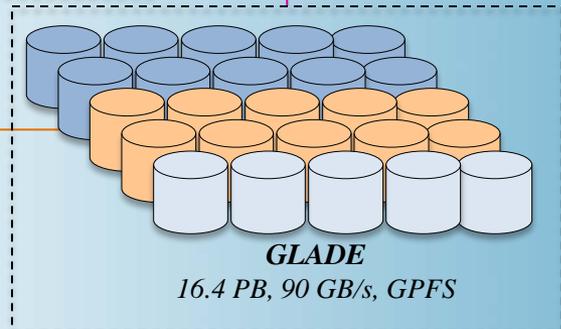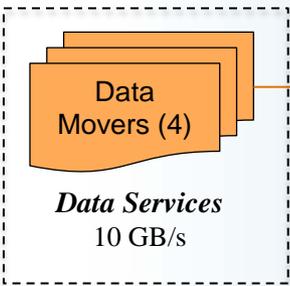
**Anke Kamrath (anke@ucar.edu)**
**Director, Computing Operations and Services , NCAR**

Computational & Information Systems Laboratory

NCAR
NATIONAL CENTER FOR ATMOSPHERIC RESEARCH

NSF

# Overview

- **Yellowstone - NCAR Data Intensive Computing Environment**

- **Globus Online**
  - Data Transfer Service
  - NCAR Cloud - Data Sharing Service

- **Large-scale Data Collections**
  - RDA – NCAR Research Data Archives
  - ESGF and CMIP – Earth System Grid Framework
  - CMIP Data Analysis Platform

- **"Herding" data from the NCAR researchers and Labs**
  - DSET – Data Stewardship Engineering Team

**Yellowstone**
*4,536 Sandy Bridge*

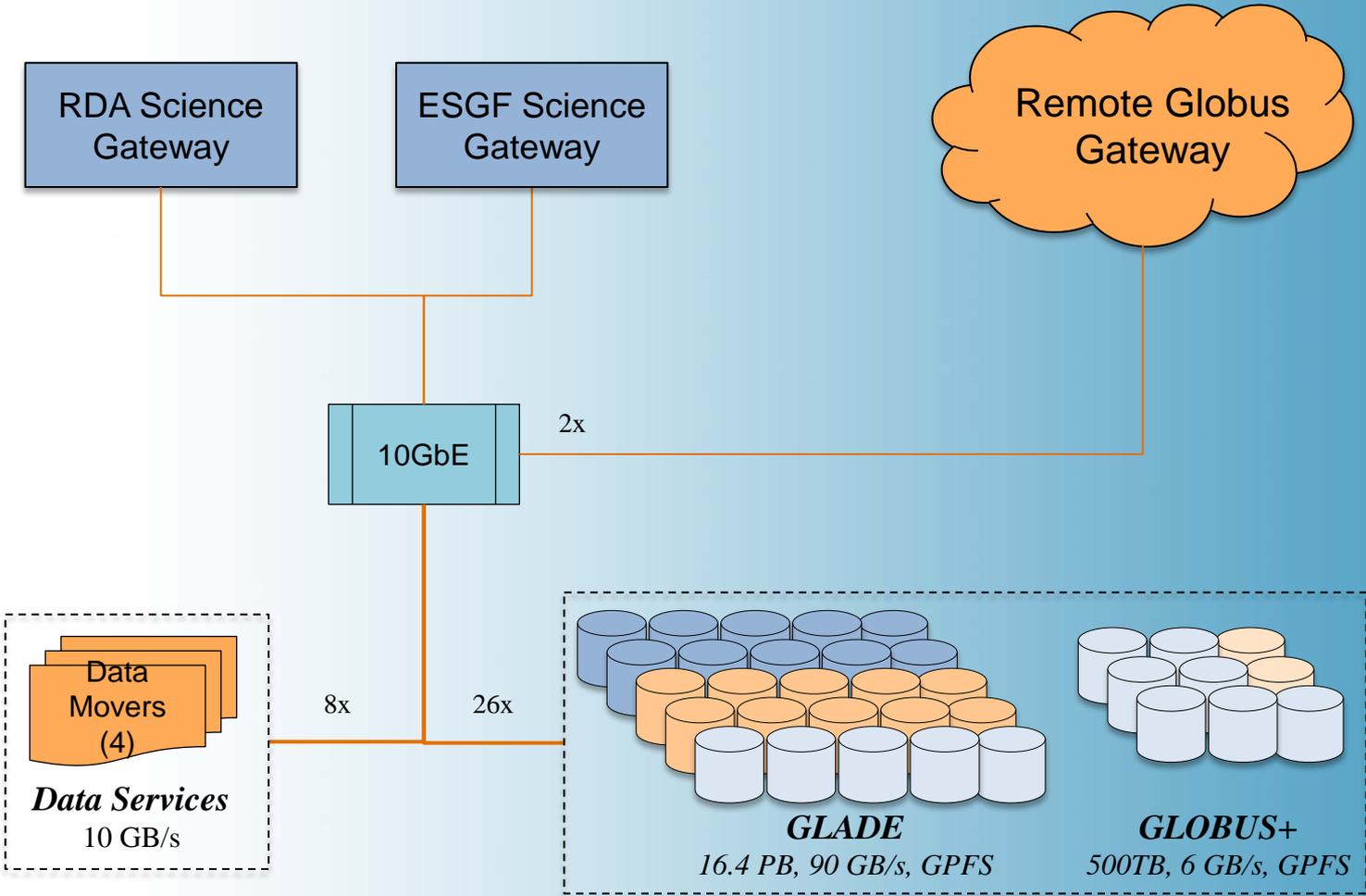72 nodes/rack — 72x — SX6036 (4x) — 72x — FDR SX6536 (9x)

Full Fat Tree IB

63x

72 nodes/rack — SX6036 (4x) — 72x

Service Nodes

Login Nodes

FDR SX6512 — 24x — FDR SX6512

26x

**Pronghorn**
*Intel KNC*
16 nodes — 16x / 16x

**Caldera**
*NVIDIA K20X*
16 nodes — 16x / 16x

**Geyser**
*1TB nodes w/NVIDIA K5000*
16 nodes — 16x / 16x

Data Movers (4)

**Data Services**
10 GB/s

8x    26x

**GLADE**
*16.4 PB, 90 GB/s, GPFS*

10GbE

HPSS

**HPSS Archive**
35PB, 12.5 GB/s

10x

2x

10GbE

2x

**Mesa Lab (Boulder)**

HPSS

**HPSS**
*Disaster Recovery*

Computational & Information Systems Laboratory

NCAR

NSF

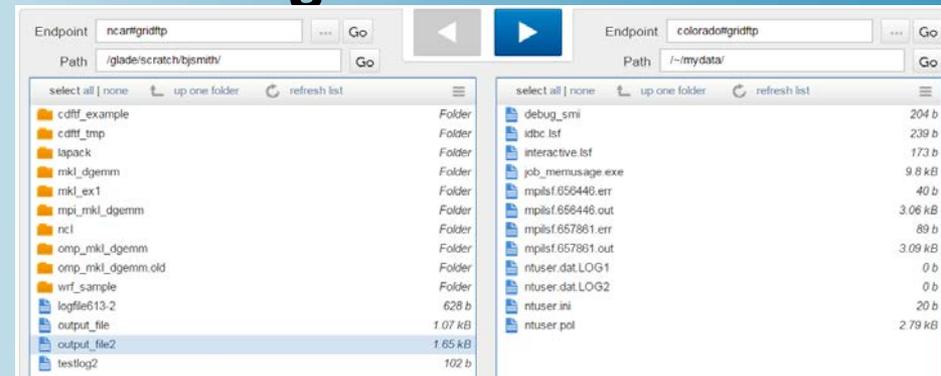NATIONAL CENTER FOR ATMOSPHERIC RESEARCH

FDR    10GbE

# Data Transfer Gateway

# *Data Transfer Services*

- **Globus Online Endpoints**
  - launch and forget data transfers
  - Access with users UCAS account and token
    - ncar#gridftp
  - Access with users XSEDE Account
    - xsede#ncar
    - xsede#glade
  - Web UI, CLI, REST API
  - Globus Connect for transfer to/from your desktop
- **gridftp, globus-url copy, scp/sft, bbcp**
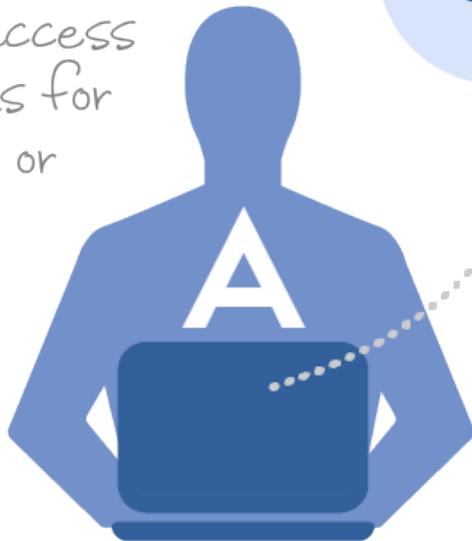- **HSI/HTAR for HPSS access through LSF**

# *Data Sharing Services*

- *ncar#datashare*
- **Globus Plus implementation**
- **data sharing allocations for self-publishing or data delivery**
- **data owner controls access**
  - can create groups for access control
  - can share 'read-only' or 'read-write'
- **user can create custom access interfaces**
  - CLI or REST API
- **500 TB of useable space**
- **> 6 GB/s bandwidth available**
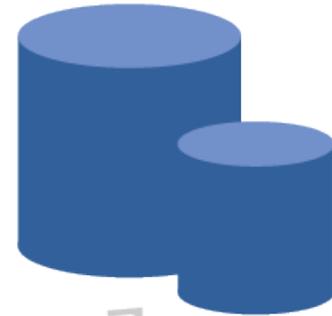- **1/2014-4/2015:  >160TB transferred.  (average filesize >3GB)**

**2** Globus tracks shared files – no need to move files to cloud storage

data source

**1** A user selects file(s) to share and sets access permissions for individuals or groups

**3** User B logs in to Globus to access shared file(s)

A

B

NCAR
NATIONAL CENTER FOR ATMOSPHERIC RESEARCH

NSF

# *Data Sharing Use Cases*

- **Delivery of data from non NCAR users for publication in a Science Gateway**
  - Periodic delivery of data sets from NCAR collaborators for publication through the ESGF Science Gateway
- **Delivery of 3D visualization to non NCAR users**
- **Publication of supporting data associated with publication**
- **Share a file or data set with a non NCAR collaborator**

# *Large NCAR Data Collections*

- **RDA – NCAR Research Data Archives**

- **ESGF and CMIP – Earth System Grid Framework**

  – CMIP Data Analysis Platform *(Coming Soon)*

# What is the RDA?

- 600+ distinct datasets for climate and weather research, 8M Files, 1.8 Pb
- Collections: ocean & atmosphere observations, analyses, reanalyses, operational NWP outputs
- Science educated staff
- Free and open access





**Unique Users**

**Data Delivered**

# *Current RDA architecture*

Computational & Information Systems Laboratory

**External Users**
**~10K/year**

- Metadata Access
- Submit Subset Requests
- Check Request Status
- Data Download
- Data Analysis

**globus online**

**RDA Web Interface and Services**

**THREDDS**
**unidata**

**RDA Web and Data Server (Migrating to VM environment in 2015)**

**RDA Metadata DB**

**MySQL**

**NCAR/CISL High Performance Computing**

**Internal Users**
**~100s/year**

- Metadata Access
- Submit Subset Requests
- Check Request Status
- Data Analysis

**CISL/NCAR Central File System 500 TB of RDA Dataset Holdings**

**CISL/NCAR HPSS Tape Library 1.8 PB of RDA Dataset Holdings**

**NCAR**
NATIONAL CENTER FOR ATMOSPHERIC RESEARCH

**NSF**

**Legend**

| | | |
|---|---|---|
| ← Data → | | |
| ← Metadata/User Interactions → | | |

11

# Private Cloud: Future RDA/NCAR architecture?

External Users ~10K/year

- Metadata Access
- Submit Subset Requests

globus online

RDA Web Interface and Services

RDA Web and Data Server (VM environment)

RDA Metadata DB

MySQL.

Internal Users ~100s/year

- Metadata Access
- Submit Subset Requests
- Check Request Status
- Data Analysis

Direct data analysis on "Private Cloud" for select external users?

NCAR/CISL High Performance Computing

Analysis

User Configured Data Analysis VM

NCAR/CISL Private Cloud?

CISL/NCAR Central File System 500 TB of RDA Dataset Holdings

CISL/NCAR HPSS Tape Library 1.8 PB of RDA Dataset Holdings

NCAR
NATIONAL CENTER FOR ATMOSPHERIC RESEARCH

Computational and Information Systems Laboratory

CISL

NSF

**Legend**

Data

Metadata/User Interactions

# RDA Usage metrics (1/2)



**Yearly Customized RDA User Access from Web Interface**
Automated Archive Subsetting and Format Conversion Requests

Legend:
- Data Volume Processed
- Data Volume Delivered
- Number of Users Served
- Number of Requests Processed

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

Legend:
- Data Volume Accessed
- Data Volume Delivered
- Number of Users Served

In 2014
- 17+ PB (virtual) processing
- 4000 users received data subsetting and format conversion services
- 45K requests were processed and 380 TB of data delivered
- 7300 users downloaded 750TB of data from archive files

# 2014 top 10 RDA datasets, ranked by number of unique users.

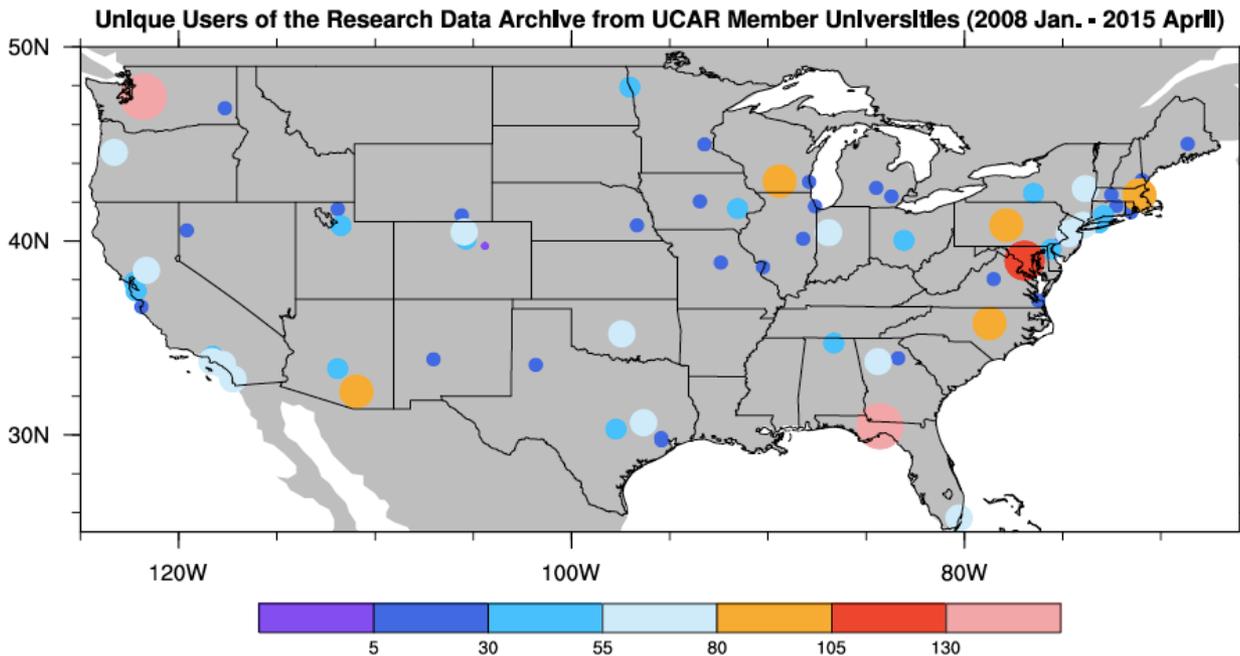| Dataset Title | Rank | Unique Users | Amount (TB) |
|---|---|---|---|
| NCEP FNL Global Tropospheric Analyses, 1999 - ongoing | 1 | 4530 | 146 |
| NCEP Climate Forecast System Global Reanalysis, 1979 - ongoing | 2 | 3156 | 505 |
| NCEP Global Upper Air and Surface Weather Obs. 1997 - ongoing | 3 | 943 | 33 |
| ECMWF-Interim Global Reanalysis, 1979 - ongoing | 4 | 490 | 137 |
| NCEP/NCAR Global Reanalysis, 1948 - ongoing | 5 | 473 | 6 |
| Japanese 55-year Global Reanalysis, 1957 - ongoing | 6 | 381 | 169 |
| NCEP North American Regional Reanalysis, 1979 - ongoing | 7 | 361 | 54 |
| Historical Unidata (IDD) Gridded Model Data, 2002 - ongoing | 8 | 348 | 29 |
| International Comprehensive Ocean-Atmosphere Data Set (ICOADS) | 9 | 268 | 1 |
| | 10 | 239 | 3 |



Unique Users of the Research Data Archive from UCAR Member Universities (2008 Jan. - 2015 April)

# *RDA DOI support*

- **Mint DOIs using EZID and managed by DataCite**
- **Management requirements set from 6 RDA use cases**
  - http://rda.ucar.edu/#!data-citation/use-cases
- **DOIs currently assigned to 49 RDA dataset collections**
- **Formal data citation in publications are supported with dashboard tool that generates citation text – see below.**

---

**For Data Accessed on 2015-03-16:**

**Dataset Citation:**  `RIS`

Gilbert P. Compo, ., et al. 2015, updated yearly. *NOAA/CIRES Twentieth Century Global Reanalysis Version 2c*. Research Data Archive at the National Center for Atmospheric Research, Computational and Information Systems Laboratory. http://dx.doi.org/10.5065/D6N877TW. Accessed 16 Mar 2015.

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

**Data Access Detail:**
9 files converted to format:

    pgrbanl_mean_1857_PRES_tropopause.grib.nc
    pgrbanl_mean_1858_PRES_tropopause.grib.nc
    pgrbanl_mean_1859_PRES_tropopause.grib.nc

# *RDA Globus integration*

- **Programmatic Globus endpoint shares**
- **Managed endpoints**
  - general archive access (Starting 4/28/15)
  - custom data orders (e.g., subset requests)
- **Access managed via Globus CLI and API**
- **Single sign-on to Globus via RDA account credentials**
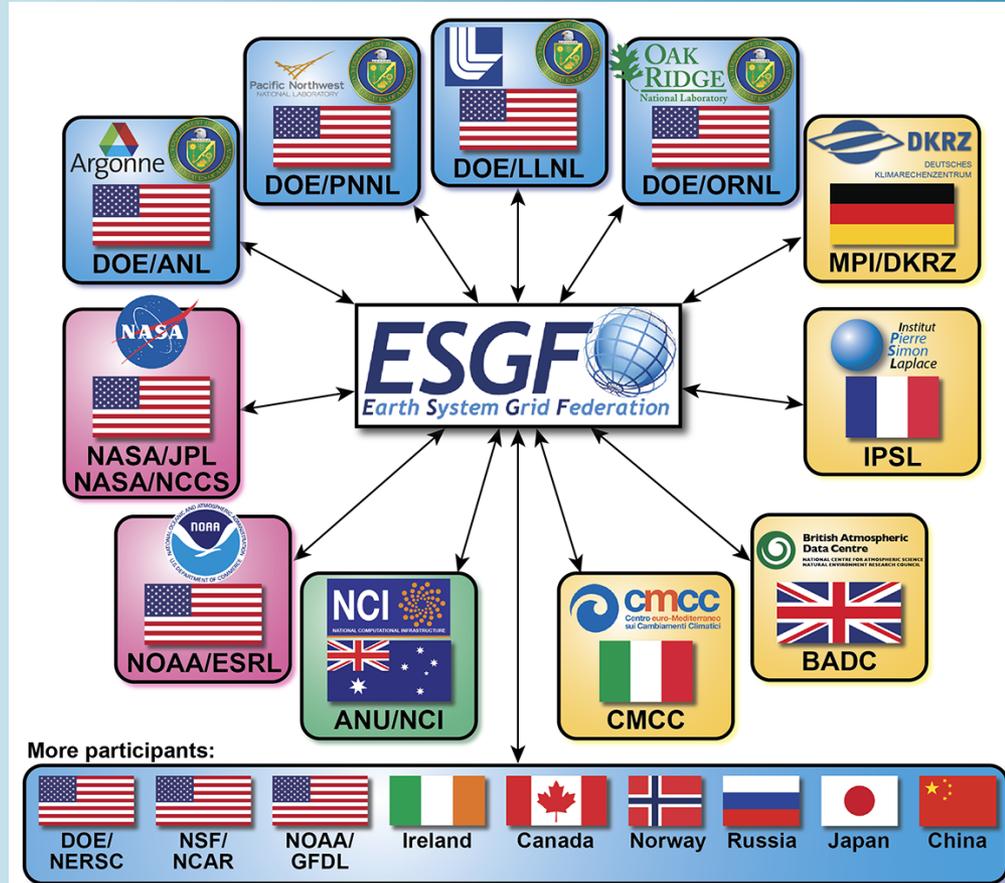
**Usage:  Nov 2014 – current, only data orders**

- **49 users**
- **14 TB transferred**

# *ESGF and CMIP*

- **ESGF – Earth System Grid Federation**
  - international collaboration for the software that powers most global climate change research, notably assessments by the Intergovernmental Panel on Climate Change (IPCC).

- **CMIP – Coupled ModeL Intercomparison**
  - protocols enable the periodic assessments carried out by the IPCC
  - Big Data!!!
    - CMIP5 – 2PB distributed across planet
    - CMIP6 (coming in ~2018) – 20PB
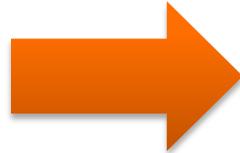
# *CMIP Analysis Platform*

- **For big data –** *"Bring compute to the data, not data to compute"*

- **Challenges:**
  - University researchers don't have infrastructure to download data and analyze
  - Takes weeks-to-months to download data.

- **Creating** *"NCAR CMIP Analysis Platform"*
  - House high-value CMIP5 (and soon CMIP6) data on GLADE
  - Also download and manage *per-request* data on GLADE
  - Advertise to US researchers, make allocations, create user accounts
  - Provide Analysis platforms (with Matlab, IDL, NCL, etc) attached to high-speed disk.

# *Working at NCAR for Comprehensive Cross-organizational Data Management*

# The
# Data Stewardship Engineering Team (DSET)

Computational & Information Systems Laboratory

CISL

NCAR
NATIONAL CENTER FOR ATMOSPHERIC RESEARCH

NSF

# The vision:

data.ncar.edu → Single front door to <u>ALL</u> data, software, data services

# *Definition*

## Data

Digital assets intended for <u>scientific community use</u>, including files and metadata, publications, reports, images, software (visualization, analysis, model codes), and related data services.

NCAR UCAR

**CLIMATE DATA** | **ANALYSIS TOOLS** | **MODEL EVALUATION** | **EXPERT CONTRIBUTORS**
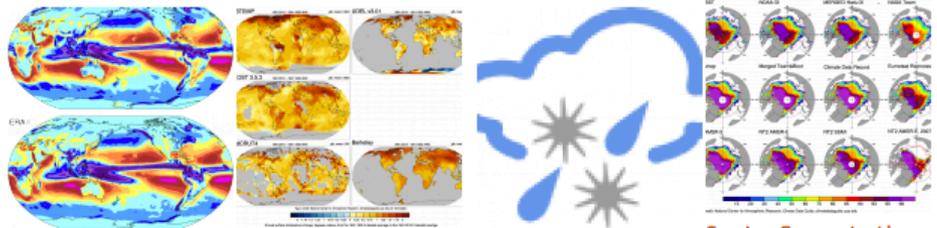
## Data Discovery Guided by Experts > >

Search and access 164 data sets covering the Atmosphere, Ocean, Land and more. Explore climate indices, reanalyses and satellite data and understand their application to climate model metrics. This is the only data portal that combines data discovery, metadata, figures and world-class expertise on the strengths, limitations and applications of climate data. Discover it now.

See data pages with guidance from these experts:

Dai, Aiguo   Norris, Joel  Willis, Josh  Meier, Walter   Dee, Dick   von Schuckmann, Karina

Banzon, Viva   Randel, Bill  Arkin, Phil  Vicente-Serrano, Sergio M.  Tian, Baijun  Kimball, John

## Data Set Overviews > >

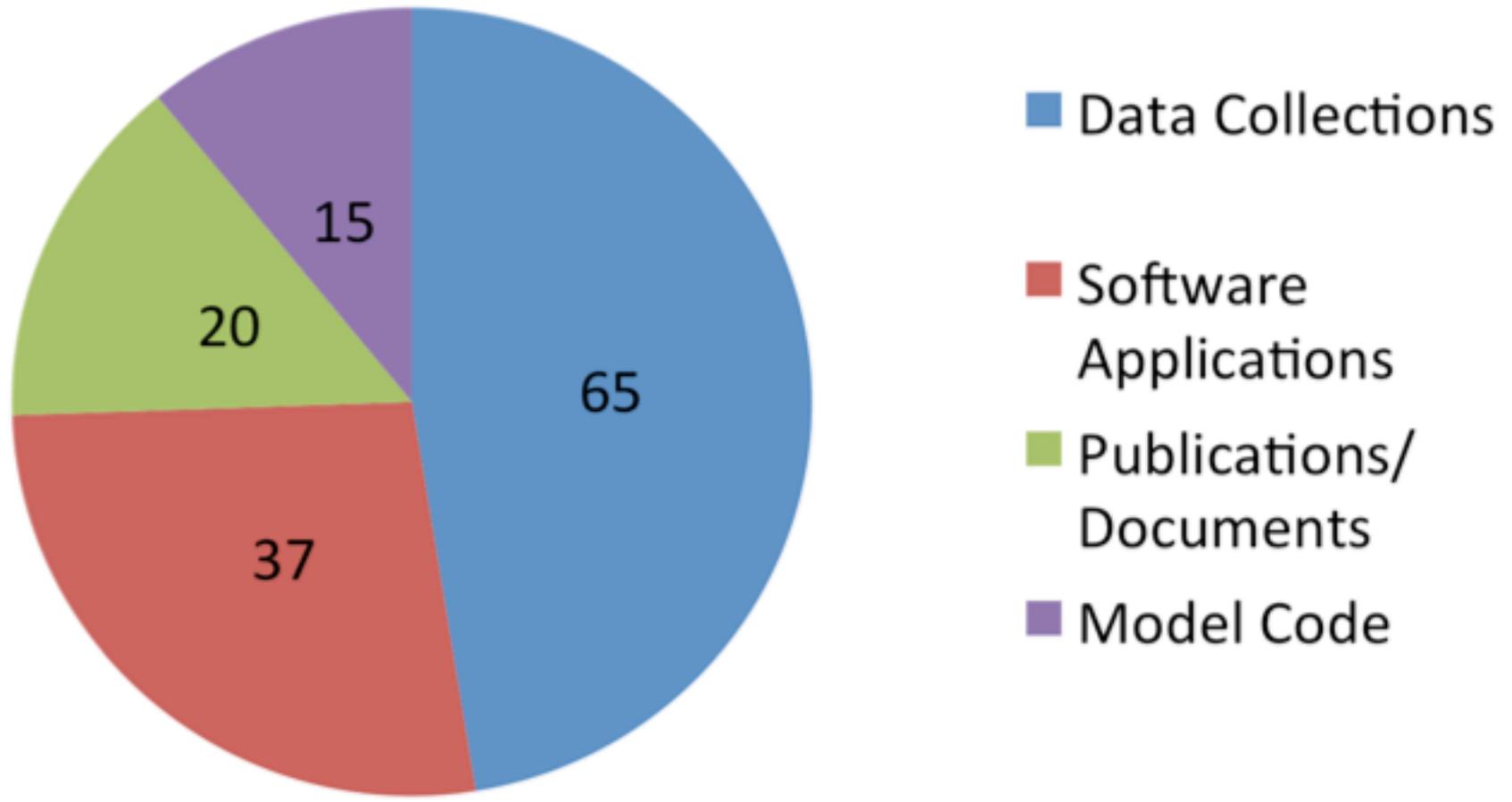Compare the attributes, strengths and limitations of multiple data sets.

Atmospheric Reanalysis: Overview & Comparison Tables

Global Temperature Data Sets: Overview & Comparison Table

Precipitation Data Sets: Overview & Comparison table

Sea Ice Concentration data: Overview, Comparison table and graphs

---

**NCAR Climate Analysis**

### Reanalysis Data

- Includes:
  - **Derived atmospheric mass, mo ECMWF and NCEP/NCAR**
  - Japanese Meteorological Agency (
  - European Center for Medium-Rang reanalyses (ERA-15 & ERA-40) an through 2010)
  - National Center for Environmenta Research reanalysis (NCEP/NCAR
  - NASA's Modern Era Retrospective
  - National Center for Environmenta (CFSR) (New!)
  - For information on all reanalyses,
  - For information on how budget pr

### Satellite Data

- Includes:
  - Adjusted ERBE and CERES retrieva
  - Data sets from a number of differe
    - ERBE
    - CERES
    - ISCCP
    - SSM/I
    - TRMM

### Surface Data

AirDat weather instruments onboard re Another focuses of such operational sys urban impacts on the weather.

In contrast to most computer weather even recognize the existence of cities ir models resolve the large–scale effects of the cities, and some even

more info

---

NCAR CGD's Climate An

NCAR UCAR · RAI

M
ACn
Wh

CIS visu CISI

SO
CIS

D
Co
A v
ge

CE
Ou
Cli

CIS
CIS
ob
su
to

Cli
NC
an
Th
fac

Cli
Se
mo
ap
wea
an

CO
Th
pu

Ea
Se
gri

airbor
and da
observ

Our co
impac
weath
data a
and im
resear

EOL Data Servi
NCAR's Earth C
archive of data
including COD

Global Change Master Dir
NASA maintains an NCAR

inte
disc

Em

www2.ucar.edu/research-resources/data-archive-

# Goals

- Inventory <u>all</u> community-based digital assets
- Assess feasibility and resource requirements
- Establish the technical and organizational foundation

# Inventory of NCAR digital collections

Collection = a group of assets managed by a system, e.g. RDA is one Data Collection



- Data Collections
- Software Applications
- Publications/ Documents
- Model Code

Data Collections: observations, images, model output
Software Applications: data processing and analysis

# *Challenges for success*

- **Level of preparedness varies greatly across NCAR**

  – Especially, regarding metadata standards

- **Need more data engineering expertise**

- **Need organization-wide digital asset consulting**

- **Need sustainable and cost effective technical solutions**

# *Stepping Up to the Challenge*

## We are:

- Reviewing, testing, and prototyping available technologies

- Making a comprehensive metadata needs assessment, by experts, across NCAR

# *Questions?*