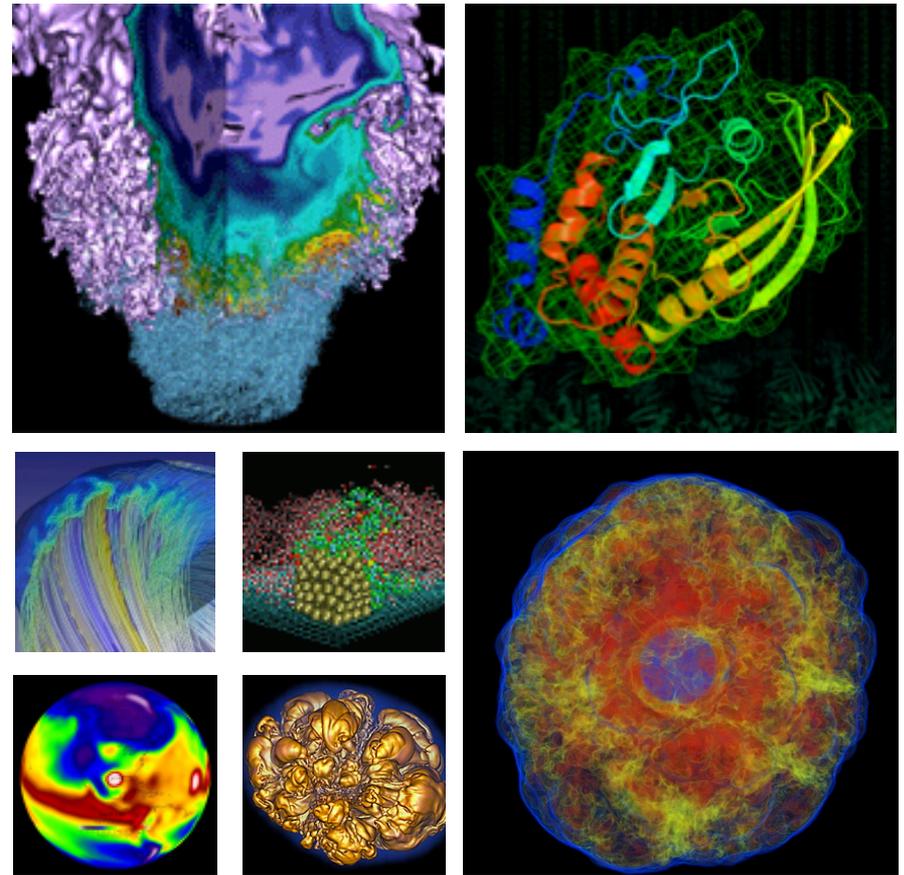


How To Share Data With Your Collaborators

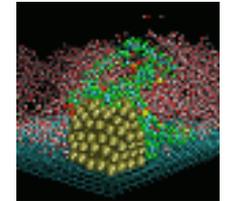
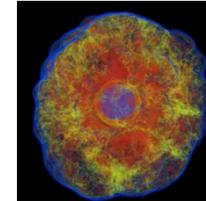
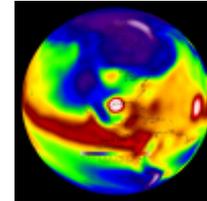
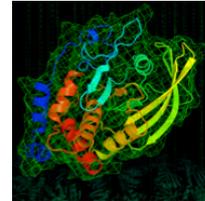
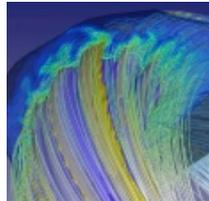
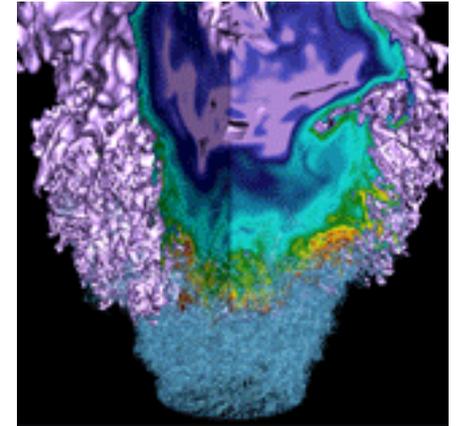


Shreyas Cholia
Data and Analytics Services, NERSC

Joint Facilities User Forum on Data-Intensive
Computing, Oakland, June 18, 2014

- **Introduction**
- **Data Sharing over HTTP**
- **APIs**
- **High Performance Data Access**
- **Data Management Systems**
- **We'll try to do a few live demos along the way**

Introduction



- **Data is still heavily siloed and inaccessible.**
 - **Data sits on a machine somewhere, and you give people local accounts to access it.**
 - **Good luck combining multiple datasets**
 - **Does not scale!**
-
- **This is 2014 – we can do better!**

- **OSTP requirements around data management**
 - All proposals submitted to the Office of Science for research funding are required to include a Data Management Plan (DMP).
 - **DMPs must provide a plan for making all research data displayed in publications resulting from the proposed research digitally accessible at the time of publication.**
 - researchers that plan to work at an Office of Science User Facility as part of the proposed research should consult the published data policy of that facility.

Data Sharing = Better Science



J.M. Wicherts, M. Bakker, and D. Molenaar:

Willingness to Share Research Data Is Related to the Strength of the Evidence and the Quality of Reporting of Statistical Results

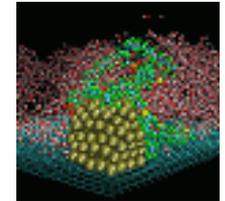
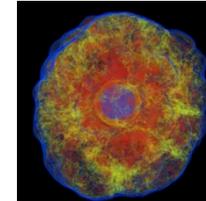
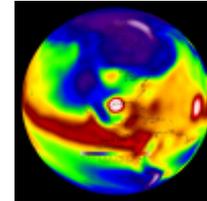
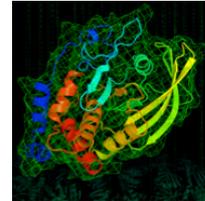
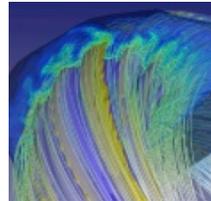
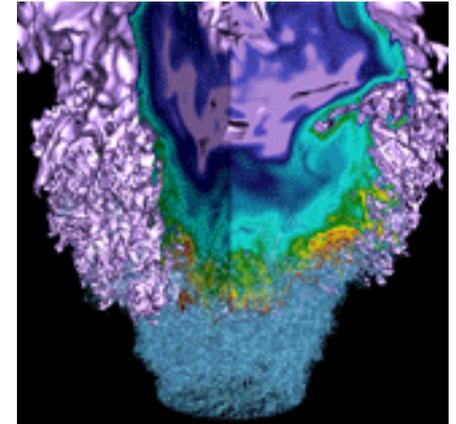
***PLoS ONE*, 6(11): e26828, 2011, doi:10.1371/journal.pone.0026828.**

Content for this slide courtesy Greg Wilson, Software Carpentry

- Common access to data enables collaboration
- Opens up the possibility of cross comparisons, combining multiple datasets from different sources
- Enables new science use cases



Data Sharing over HTTP

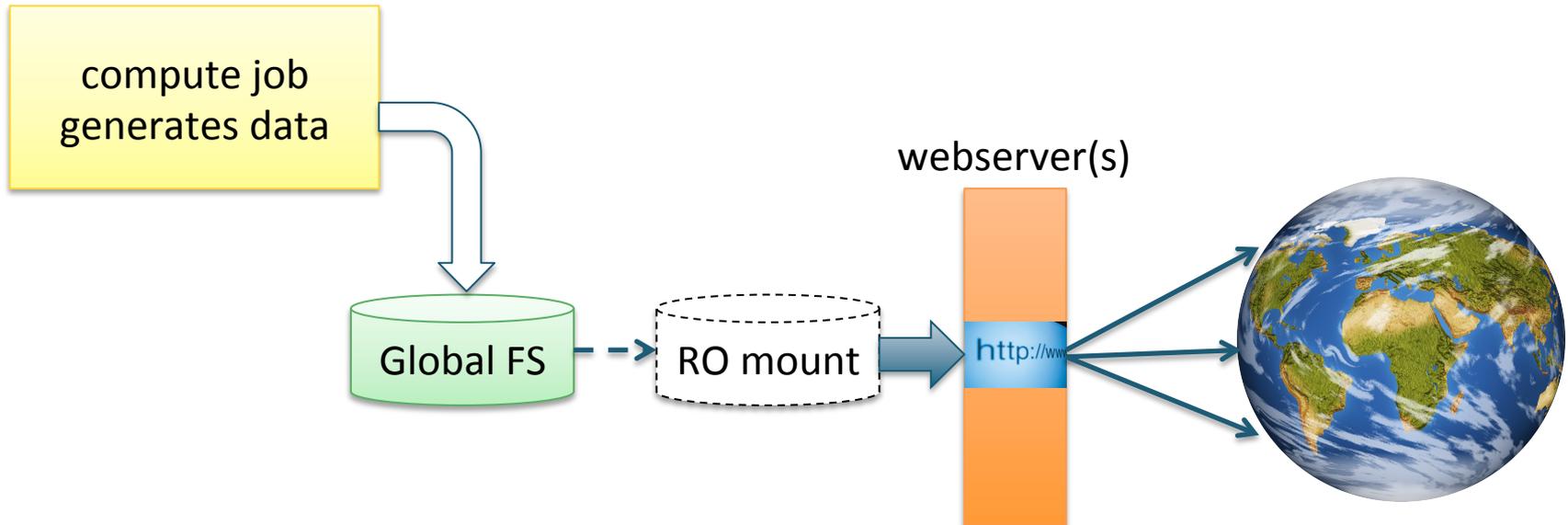


Sharing Data Over The Web



- **Provide projects with a landing spot for data within the center boundary**
 - High performance access to filesystems, backend resources etc.
- **Serve up landing spot via web server**
 - Minimally provides a very simple publication sharing system

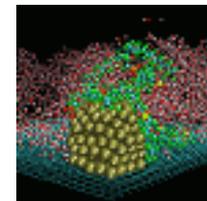
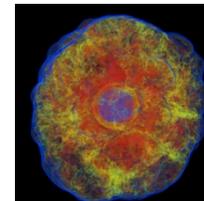
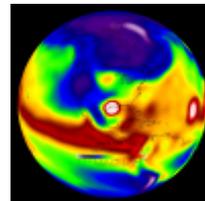
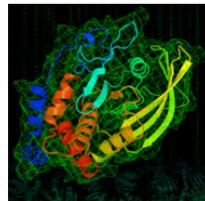
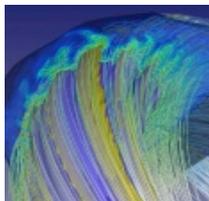
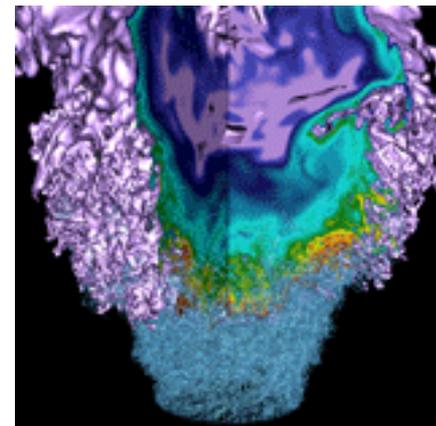
Data Sharing Over The Web



- **Store data on persistent global filesystem**
- **Designate an area for export**
- **Mount export area RO on a web node using a shared FS**
- **Run a webserver on web node to share with world**
- **Add AuthNZ to web server configuration**
- **Demo time!**

- **Authorization: Access privileges for a given identity**
 - What is the user allowed to do?
- **Data portals see some combination of the following sharing modes:**
 - Public – data are publicly viewable to all users
 - Protected – data are only visible to members within a designated group or collaboration
 - Private – data are private to the individual user

APIs



APIs Drive Programmatic Access to Data

The NERSC logo is a dark blue rectangle with a white starburst effect behind the text "NERSC" in white, bold, sans-serif font.

- **Make your data available as a Web API to enable more fine-grained queries**
- **Use HTTP Method + URLS + parameters to query and access data**
- **Users just download what they need, server does heavy lifting in terms of search and sub-selection**
 - eg. Google maps - you don't download the map just the tiles you are interested in.
- **Enable rich interactions and applications around your data**
- **eg. OpenDAP connects NETCDF/HDF5 datasets with web clients**
 - Demo: <http://portal.nersc.gov/pydap>

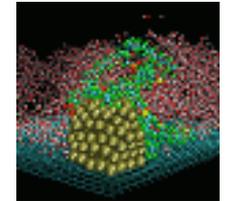
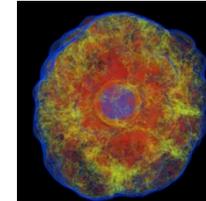
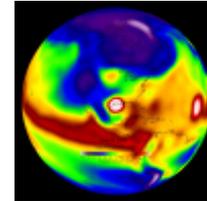
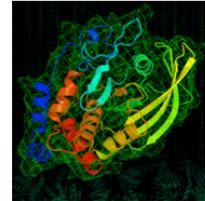
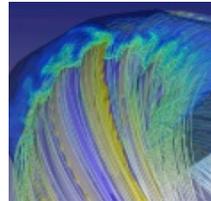
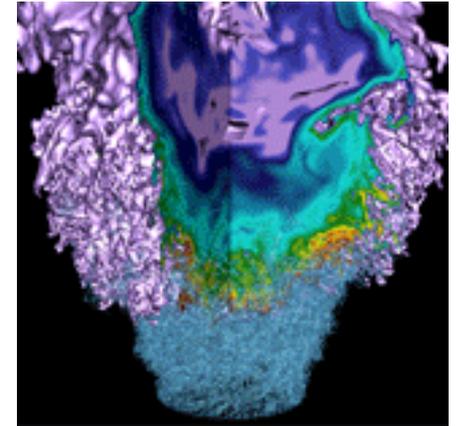
- **HTML5/JavaScript + Server Side APIs to build data browsers and viewers**
- **Make callouts to REST API for all server-side information**
- **eg. <https://openmsi.nersc.gov>**

Getting started – Hello World API

> pip install flask

```
flask_demo.py
1  from flask import Flask
2  from flask import jsonify
3  app = Flask(__name__)
4
5  # http://localhost:5000/object_id - returns object
6  @app.route('/<object_id>')
7  def get_data(object_id):
8      obj = get_object_from_db(object_id)
9      return jsonify(obj)
10
11 # http://localhost:5000/ - returns all objects
12 @app.route('/')
13 def get_all_data():
14     obj_list = get_data_from_db()
15     return jsonify(obj_list)
16
17 if __name__ == "__main__":
18     app.run()
19
```

High Performance Data Sharing



High Performance Data Sharing

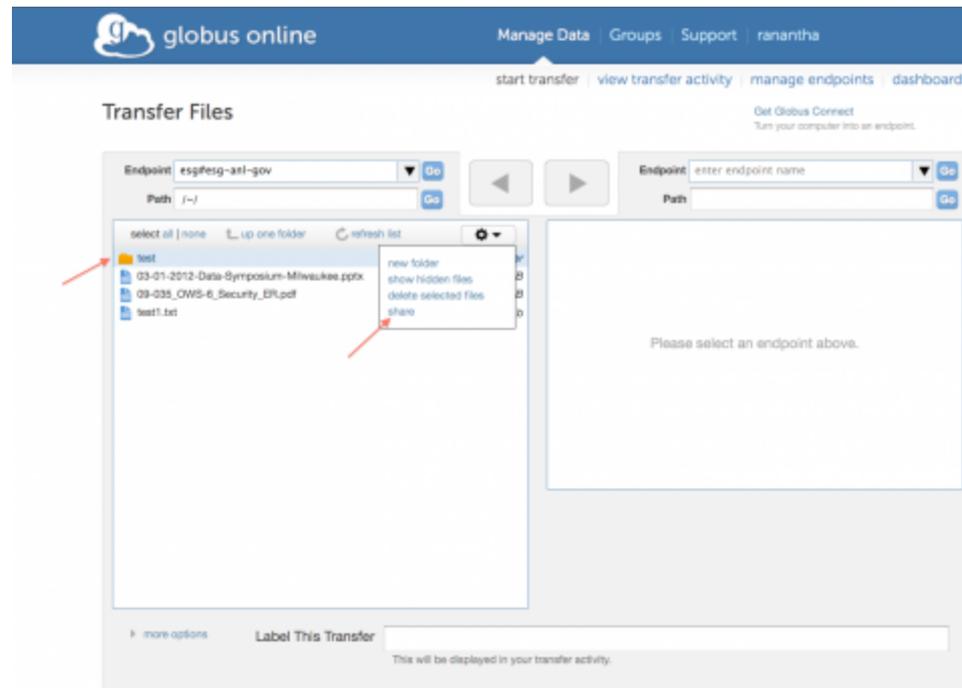


- **Globus**
- **webdav + xrootd**
- **Aspera (Commercial Tool)**

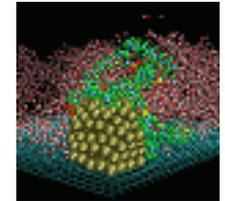
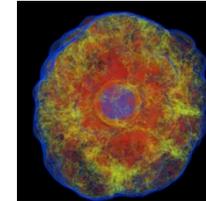
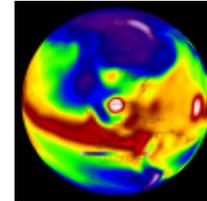
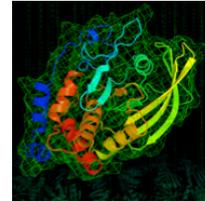
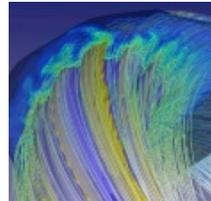
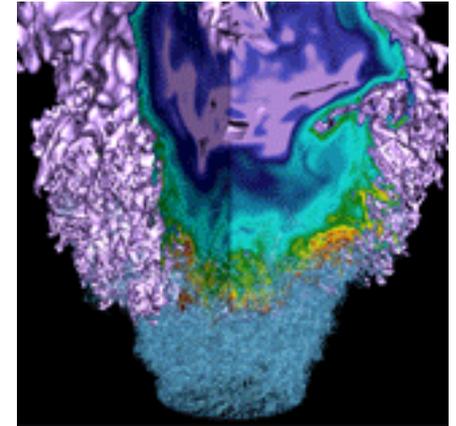
Globus Sharing



- Enable high performance GridFTP access to your data
- Share with designated groups or individuals
- <https://www.globus.org/data-sharing>



Data Management Portals

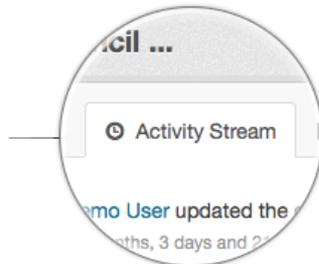
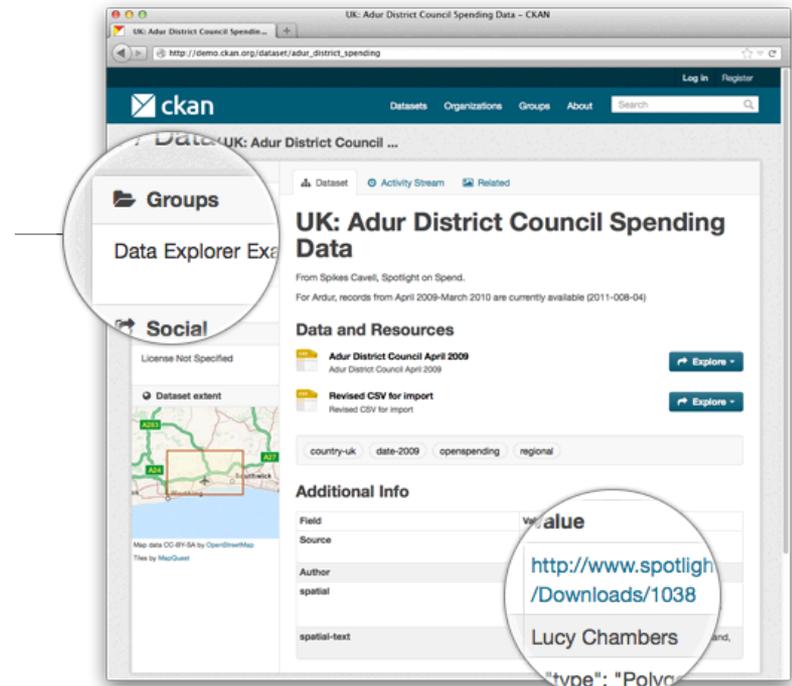
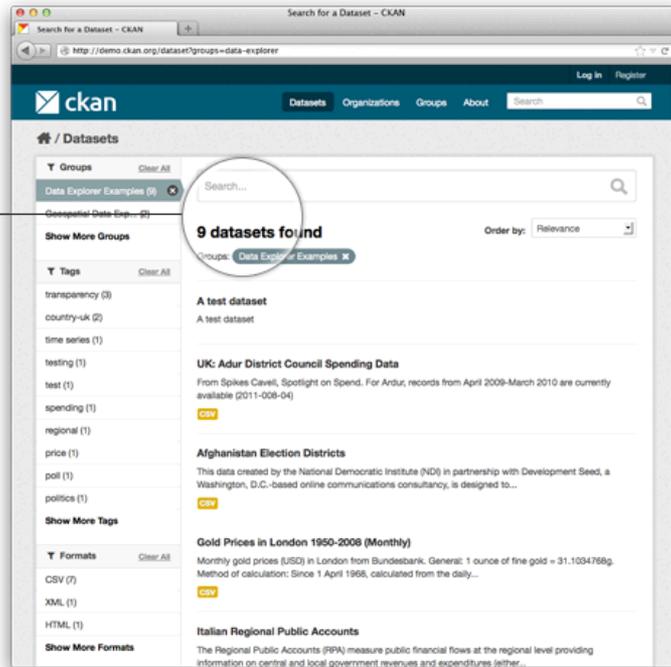


- **Data Management System (like a CMS but with a focus on Data)**
- **CKAN – open source tool for data publishing & management used by data.gov.uk and publicdata.eu**
 - Publish and Manage Data
 - Search and Discovery
 - Metadata
 - Geospatial and Visualization Features
 - Community
 - Link to Filesets or store within portal
 - Data Provenance
 - Extensibility and APIs
 - Federation

Search...

9 datasets found

Groups: Data Explorer



CKAN demo



- <http://ckan.org>

- In some cases data may be distributed across multiple locations
- There are federated portals like ESGF which enable searching across data from multiple locations
- Metadata is federated but data is ultimately downloaded from source portal

ESGF – Federated Search



Current Selections

- (x) [query:cmor](#)

Search Categories

Project

Institute

Model

SubModel

Instrument

Experiment Family

Experiment

SubExperiment

Time Frequency

Product

Realm

Variable

Variable Long Name

Search

Examples: *temperature*, "surface temperature", *climate AND project:CMIP5 AND variable:hus*.
To download data: add datasets to your Data Cart, then click on *Expand* or *wget*.

[Temporal Search](#)
[Clear search](#)
[constraints and datacart](#)
[Search Help](#)
[Search Controlled](#)
[Vocabulary](#)

Search All Sites Show All Replicas Show All Versions

< 1 2 3 ... 6518 6519 > displaying 1 to 10 of 65186 search results

Display datasets per page

[Add All Displayed to Datacart](#) [Remove All Displayed from Datacart](#)

Results

Data Cart

[project=CMIP5,model=EC-EARTH consortium,experiment=10- or 30-year run initialized in year 1960,time_frequency=day,modeling_realm=sealce,ensemble=r10i3p1,version=20140318](#)

Data Node: [esg-dn1.nsc.liu.se](#)

Version: 20140318

Description: EC-EARTH model output prepared for CMIP5 10- or 30-year run initialized in year 1960

Further options: [Add To Cart](#) [Visualize and Analyze](#) [Model Documentation](#)

[project=CMIP5,model=EC-EARTH consortium,experiment=10- or 30-year run initialized in year 1960,time_frequency=day,modeling_realm=sealce,ensemble=r1i3p1,version=20140318](#)

Data Node: [esg-dn1.nsc.liu.se](#)

Version: 20140318

Description: EC-EARTH model output prepared for CMIP5 10- or 30-year run initialized in year 1960

Further options: [Add To Cart](#) [Visualize and Analyze](#) [Model Documentation](#)

The End



- **Contact:**
 - Shreyas Cholia: scholia@lbl.gov
- **Questions? Comments?**