

# Research Data Management

Simple Ways to Make your Research Life Easier

Tom Morrell

BE/Bi 103

October 10, 2018

# Current Research Data Practices



Most researchers store data on local computer hard drives

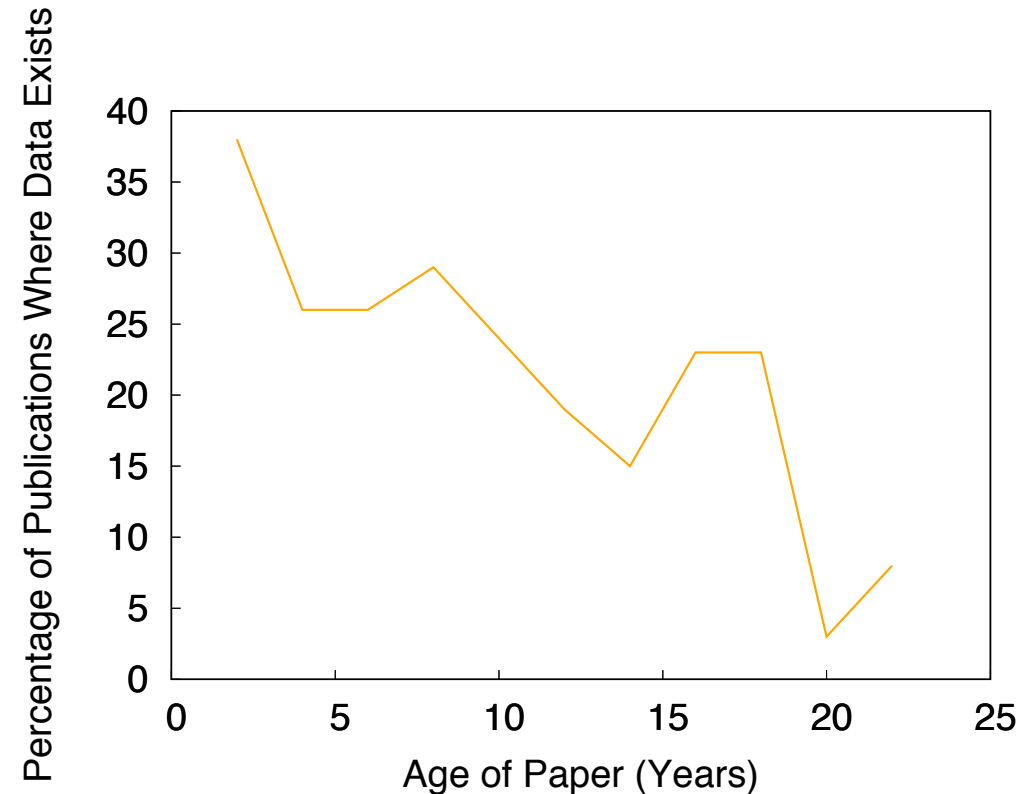
Researchers report that finding data is their biggest challenge

Akers, K. G. & Doty, J. Disciplinary differences in faculty research data management practices and perspectives. *Int. J. Digit. Curation* **8**, 5–26 (2013). (Emory)

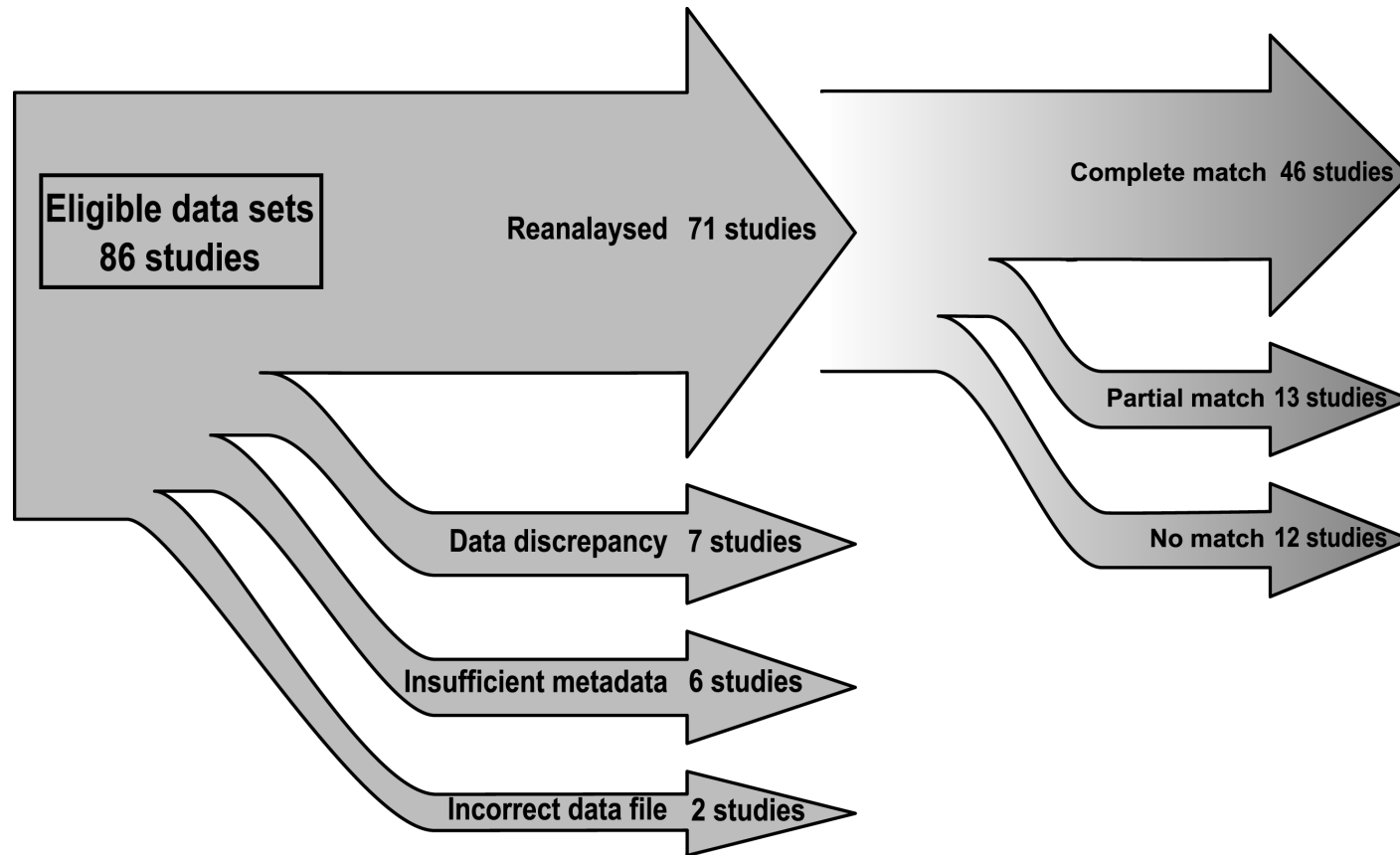
Shen, Y. Strategic Planning for a Data-Driven, Shared-Access Research Enterprise: Virginia Tech Research Data Assessment and Landscape Study. *Coll. Res. Libr.* **77**, 500–519 (2016).

# How Reusable is Research Data Today?

- Morphological characteristics of plants and animals
  - 516 publications using a specific analysis technique between 1991 and 2011
    - 25% of emails didn't work
    - 38% didn't respond to email
    - 13% didn't have data
    - 4% didn't want to share
    - Received 19% of data
    - Availability decreased with time



# Data Quality

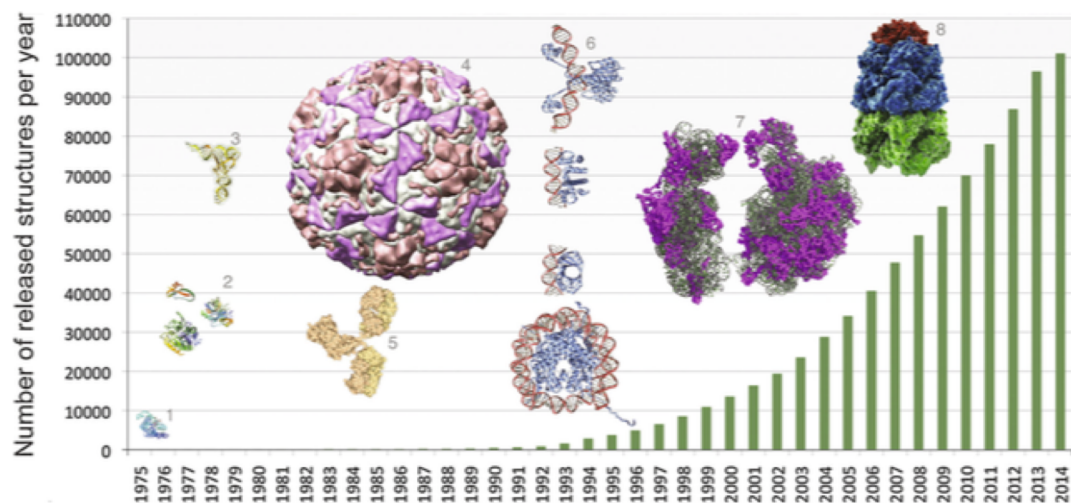


On Average, 13%  
of Papers Had  
Usable Data

# Why is it better to have data available?



[www.rcsb.org](http://www.rcsb.org)



Biological assembly 1 assigned by authors

Select a Viewer

JSmol (JavaScript)

Stoichiometry: 1

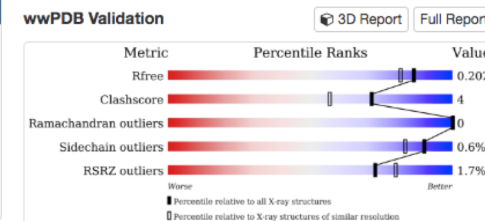
Select Orientation: Front

Select Display Mode: Secondary Structure, Subunit, Symmetry

Display Options: Style (Cartoon), Color (Secondary Structure), Surface (None)

Scripting Options: Ligands, Domain, Modification

Ligand ID	View Interactions	Ligand Electron Density	Image	Name / Formula / Weight
MTA	<a href="#">View SHJM - MTA Pocket Interaction</a>	A:401		5'-DEOXY-5'-METHYLTHIOADENOSINE C11 H15 N5 O3 S 297.334



# Why is it better to have data available?

“Digitally formatted scientific data resulting from unclassified research supported wholly or in part by Federal funding should be stored and publicly accessible to search, retrieve, and analyze.”

2013 OSTP Memo

## Data Management Plans

- Expected Data
- Data Formats and Metadata
- Access to Data
- Data Archiving

# Why is it better to have data available?

Journals require data availability:



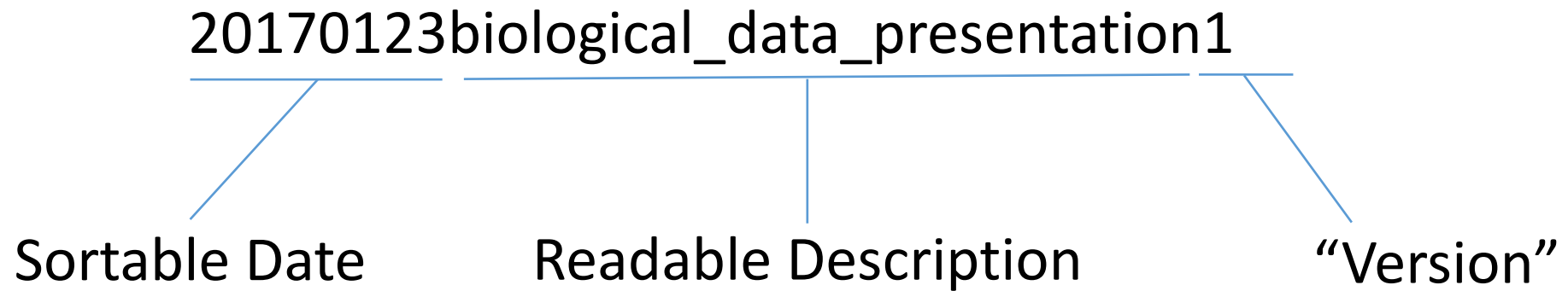
## Simple Solutions

- Choose a file naming/organization scheme
- Save reasonable files
- Use reliable storage
- Plan for sharing



# Naming

- Trying to recreate your work months/years later is hard
- Choosing a consistent naming system makes things easier



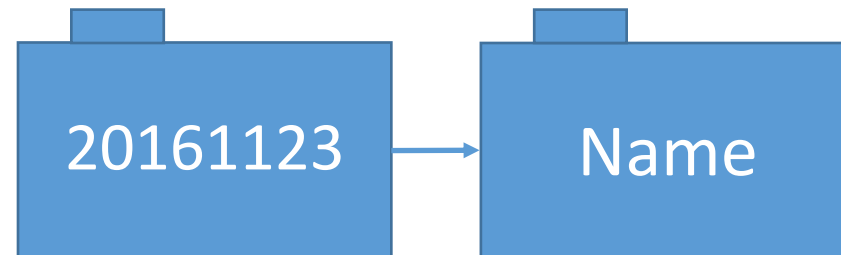
# Data Architectures

Simple

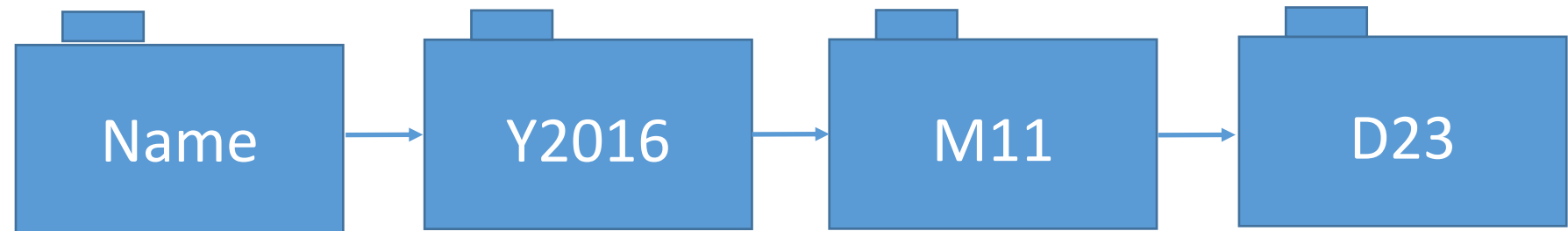


**Dataset**  
Automatically Manage Metadata/Documents  
<https://github.com/caltechlibrary/dataset>

Date Based



Complex



# Save Reasonable Files

- Human-readable text files are best (.txt, .csv)
- Non-proprietary files are better than proprietary
- Do analysis with scripts if possible
- Save both input and output files as space allows

# Active Data Storage

- Small amounts of data (GB) are easy
- TB-scale data require planning
  - Need a system that will be reliable
    - Network-Attached Storage (Local RAID array)
    - Cloud Storage

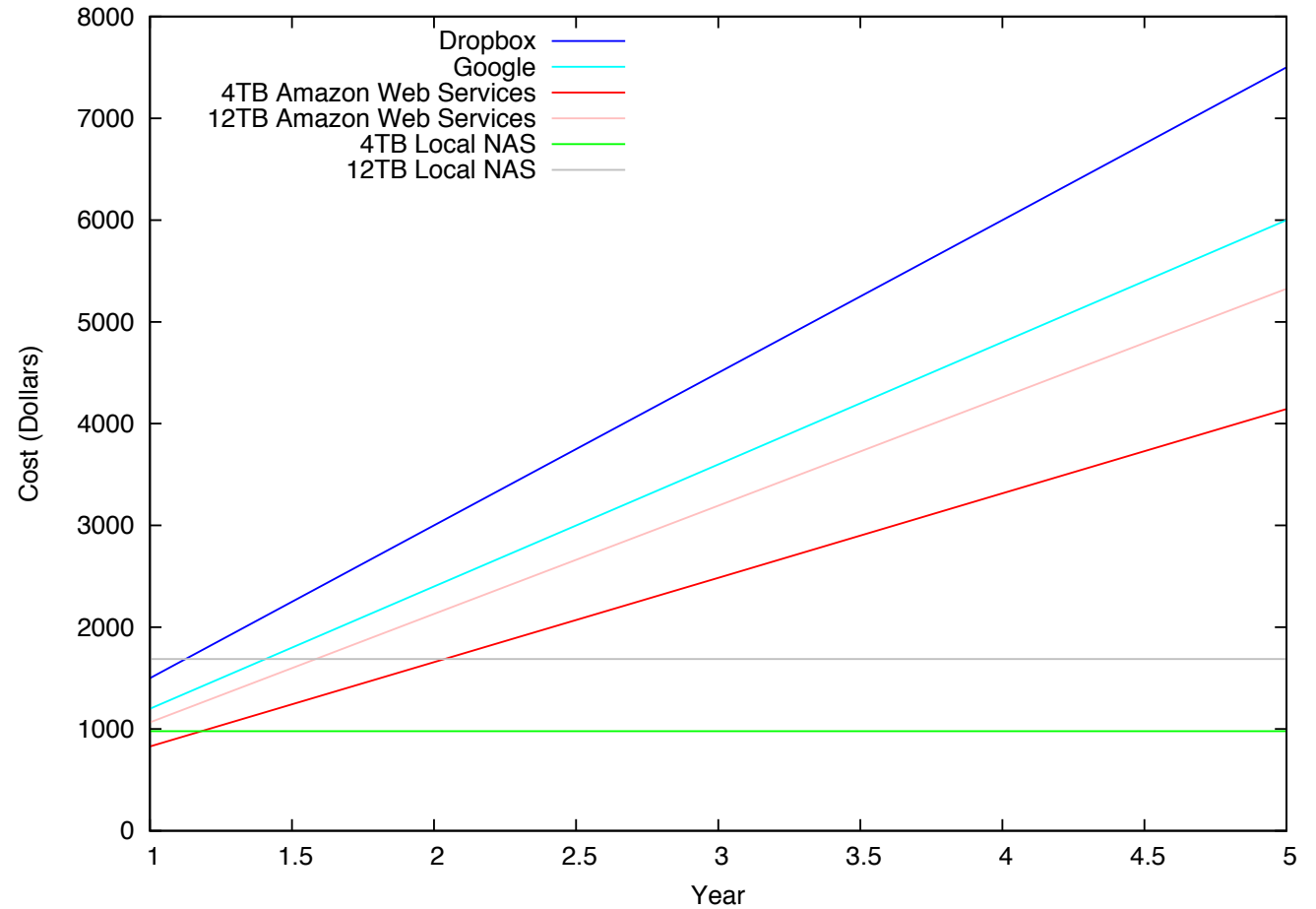
# Network-Attached Storage

- Small computer with array of hard disks
- Consumer/Prosumer devices
- Low Cost (4 TB-\$425; 42 TB-\$3000)
- Need to plan space requirements
- Need to manage



# Cloud Storage

- Defined or flexible storage
- Vendor Managed
- Continuous cost
- Limited by bandwidth
- Dependent on vendor



# Disaster Recovery

- What Happens in a Disaster?
- Use 2 mirrored NAS units in 2 locations
- Mirror NAS to cloud storage  
(Box.com - [imss.caltech.edu/box](https://imss.caltech.edu/box))



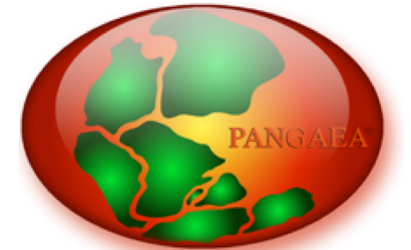
# Data Sharing

- FAIR (Findability, Accessibility, Interoperability, Reusability)
  - Subject Repositories
  - General Repositories
  - Institutional Repositories



# Subject Repositories

- Protein Data Bank
- GenBank
- Wormbase
- Pangaea
- Long Term Ecological Research Data Portal
- Good listing: [journals.plos.org/plosone/s/data-availability](http://journals.plos.org/plosone/s/data-availability)
- Thousands more: [www.re3data.org](http://www.re3data.org)



# General Repositories

The Zenodo logo consists of the word "zenodo" in a white, lowercase, sans-serif font, centered within a solid blue rectangular background.

- Zenodo (CERN-Free)
- Dryad (Nonprofit-\$120 per submission + Space)
- Figshare (20GB Max)
- Mendeley Data (Elsevier-Free)
- Dataverse (Harvard-Free)

# CaltechDATA



California Institute of Technology

Research Data Repository

- Available at [data.caltech.edu](https://data.caltech.edu)
- Easy to describe and upload files
- All records get a DOI (permanent, registered link)
- Integration with Github
- API for accessing data
- Library takes care of preserving and maintaining access to files



**GitHub**

# Discoverability

- CaltechDATA site search
- DOIs appear in DataCite search
- and Search Engines

The image shows two overlapping search results for the paper "Identifying and Quantifying Mineral Abundance through VSWIR Microimaging Spectroscopy: A Comparison to XRD and SEM" by Leask, Ellen K. and Ehlmann, Bethany L.

**Google Search Results:**

- Search query: VSWIR microimaging
- Results: About 903 results (0.37 seconds)
- Top result: <https://ssed.gsfc.nasa.gov/IPM/PDF/1046.pdf> - MICROIMAGING VSWIR SPECTROSCOPY INSTRUMENTS FOR PLANETARY EXPLORATION: MEASURING IN-SITU MINERALOGY, ICES, ORGANICS AND ...
- Second result: Using VSWIR Microimaging Spectroscopy to Explore the ... authors.library.caltech.edu/70145/ by AA Fraeman - 2016 - Related articles
- Third result: Figure 1. VSWIR Microimaging spectroscopy of a fragment of Allende ... www.hou.usra.edu/meetings/ipm2016/pdf/4097.pdf - VSWIR Microimaging spectroscopy of a fragment of Allende demonstrates the ability to map mineralogical variability in fine-grained dark materials characteristic ...
- Fourth result: IDENTIFYING AND QUANTIFYING MINERAL ABUNDANCE ... www.hou.usra.edu/meetings/ipm2016/pdf/4022.pdf - IDENTIFYING AND QUANTIFYING MINERAL ABUNDANCE THROUGH VSWIR MICROIMAGING. SPECTROSCOPY: A COMPARISON TO XRD AND SEM.
- Fifth result: VSWIR Microimaging Spectroscopy for Geologic History and ... adsabs.harvard.edu/abs/2016LPICo1980.4097E by BL Ehlmann - 2016
- Sixth result: Identifying and Quantifying Mineral Abundance Through VSWIR ... adsabs.harvard.edu/abs/2016LPICo1980.4022L by EK Leask - 2016 - Related articles
- Seventh result: Identifying and Quantifying Mineral Abundance through VSWIR ... - DOIs https://doi.org/10.22002/D1.222

**DataCite Search Results:**

- Search query: VSWIR microimaging
- 1 result
- Sorted by: Relevance
- Resource Type: Dataset
- Author: Ehlmann, Bethany L.; Leask, Ellen K.
- Title: Identifying and Quantifying Mineral Abundance through VSWIR Microimaging Spectroscopy: A Comparison to XRD and SEM
- DOI: https://doi.org/10.22002/D1.222

**CaltechDATA Site Search Results:**

- Search query: VSWIR microimaging
- 1 result
- Sorted by: Relevance
- Resource Type: Dataset
- Author: Ehlmann, Bethany L.; Leask, Ellen K.
- Title: Identifying and Quantifying Mineral Abundance through VSWIR Microimaging Spectroscopy: A Comparison to XRD and SEM
- DOI: https://doi.org/10.22002/D1.222

**Callout Box:**

Identifying and Quantifying Mineral Abundance through VSWIR ... - DOIs  
<https://doi.org/10.22002/D1.222>  
 Mar 13, 2017 - Identifying and Quantifying Mineral Abundance through VSWIR Microimaging Spectroscopy: A Comparison to XRD and SEM. Dataset.

# Citations

**TCCON data from Caltech (US), Release GGG2014.R1**

Dataset  
2017-09-08  
CaltechDATA

Download Edit

**Details**

**Authors**  
Wennberg, P. O. California Institute of Technology, Pasadena, CA (US) 0000-0002-6126-3854 ORCID  
Wunch, D. California Institute of Technology, Pasadena, CA (US) 0000-0002-4924-0377 ORCID  
Roehl, C. M. California Institute of Technology, Pasadena, CA (US) 0000-0001-5383-8462 ORCID  
Blavier, J.-F. Jet Propulsion Laboratory, Pasadena, CA (US) 0000-0002-1800-9316 ORCID  
Toon, G. C. Jet Propulsion Laboratory, Pasadena, CA (US)  
Allen, N. T. Harvard University, Cambridge, MA (US) 0000-0002-7528-8606 ORCID

**Contributors**  
HostingInstitution California Institute of Technology, Pasadena, CA (US)  
DataCurator Roehl, C. M. California Institute of Technology, Pasadena, CA (US) 0000-0001-5383-8462 ORCID  
ContactPerson Paul Wennberg wennberg@gss.caltech.edu

**Description**  
Abstract:  
The Total Carbon Column Observing Network (TCCON) is a network of ground-based Fourier Transform Spectrometers that record direct solar absorption spectra of the atmosphere in the near-infrared. From these spectra, accurate and precise column-averaged abundances of atmospheric constituents including CO<sub>2</sub>, CH<sub>4</sub>, N<sub>2</sub>O, HF, CO, H<sub>2</sub>O, and HDO, are retrieved. This data set contains observations from the TCCON station at the California Institute of Technology, Pasadena, U.S.A.

**Publication Date**  
2017-09-08

**Subjects**  
atmospheric trace gases, CO<sub>2</sub>, CH<sub>4</sub>, CO, N<sub>2</sub>O, column-averaged dry-air mole fractions, remote sensing, FTIR spectroscopy, TCCON

**DOI**  
10.14291/tcon.ggg2014.pasadena01.R1/1182415

**Version**  
GGG2014.R1

**Format**  
application/x-netcdf



<https://doi.org/10.14291/tcon.ggg2014.pasadena01.R1/1182415>

Update Record

Related Identifier(s)

- IsDocumentedBy (URL): [https://tcon-wiki.caltech.edu/Network\\_Policy/Data\\_Use\\_Policy/Data\\_Description](https://tcon-wiki.caltech.edu/Network_Policy/Data_Use_Policy/Data_Description)
- IsDocumentedBy (URL): <https://tcon-wiki.caltech.edu/Sites>
- IsPartOf (URL): <http://tcondata.org>
- IsDocumentedBy (DOI): [10.14291/tcon.ggg2014.documentation.R0/1221662](https://doi.org/10.14291/tcon.ggg2014.documentation.R0/1221662)
- IsCitedBy (DOI): [10.5194/amt-9-683-2016](https://doi.org/10.5194/amt-9-683-2016)
- IsCitedBy (DOI): [10.5194/amt-9-227-2016](https://doi.org/10.5194/amt-9-227-2016)
- IsCitedBy (DOI): [10.5194/amt-9-3491-2016](https://doi.org/10.5194/amt-9-3491-2016)
- IsCitedBy (DOI): [10.5194/amt-9-3527-2016](https://doi.org/10.5194/amt-9-3527-2016)
- IsNewVersionOf (DOI): [10.14291/tcon.ggg2014.pasadena01.R0/1149162](https://doi.org/10.14291/tcon.ggg2014.pasadena01.R0/1149162)
- IsPartOf (DOI): [10.14291/TCCON\\_GGG2014](https://doi.org/10.14291/TCCON_GGG2014)
- IsCitedBy (DOI): [10.3390/rs8050414](https://doi.org/10.3390/rs8050414)

remote sensing

Title / Keyword:  Journal: Remote Sensing

Author / Affiliation:  Section: all

Article Type: all Special Issue:

Advanced Search

Volume 8, Issue 5

Article Versions

- Abstract
- Full-Text PDF [2676 KB]
- Full-Text HTML
- Full-Text XML
- Full-Text Epub
- Article Versions Notes
- Supplementary material

Related Info

- Gonnie Schlar

Remote Sens. 2016, 8(5), 414; doi:[10.3390/rs8050414](https://doi.org/10.3390/rs8050414)

Article

**Comparison of XH<sub>2</sub>O Retrieved from GOSAT Short-Wavelength Infrared Spectra with Observations from the TCCON Network**

Eric Dupuy <sup>1,\*</sup>, Isamu Morino <sup>1</sup>, Nicholas M. Deutscher <sup>2,3</sup>, Yukio Yoshida <sup>1</sup>, Osamu Uchino <sup>1</sup>, Brian J. Connor <sup>4</sup>, Martine De Mazière <sup>5</sup>, David W. T. Griffith <sup>2</sup>, Frank Hase <sup>6</sup>, Pauli Heikkinen <sup>7</sup>, Patrick W. Hillyard <sup>8,9</sup>, Laura T. Iraci <sup>8</sup>, Shuji

48. Wennberg, P.O.; Wunch, D.; Roehl, C.; Blavier, J.F.; Toon, G.C.; Allen, N. *TCCON Data from California Institute of Technology, Pasadena, California, USA, Release GGG2014R1*; Carbon Dioxide Information Analysis Center; Oak Ridge National Laboratory: Oak Ridge, TN, USA, 2014. [Google Scholar] [CrossRef]

<https://doi.org/10.3390/rs8050414>

Citation

Email Alert

California Institute of Technology  
Research Data Repository

Dear Paul Wennberg,

Your CaltechDATA work "TCCON data from Caltech (US), Release GGG2014.R1" has been cited in:

1. Dupuy E, Morino I, Deutscher N, et al. Comparison of XH<sub>2</sub>O Retrieved from GOSAT Short-Wavelength Infrared Spectra with Observations from the TCCON Network. Remote Sensing. 2016;8(5):414. doi:10.3390/rs8050414.

This link has been added to your CaltechDATA record at [10.14291/tcon.ggg2014.pasadena01.R1/1182415](https://doi.org/10.14291/tcon.ggg2014.pasadena01.R1/1182415).

Best,  
CaltechDATA Alerting Service

Is this incorrect? Let us know at [data@caltech.edu](mailto:data@caltech.edu)

This email was sent by the Caltech Library, 1200 East California Blvd., MC 1-43, Pasadena, CA 91125, USA

[unsubscribe](#)



California Institute of Technology

# Research Data Repository

Demo

# Use Cases - Theses

Upload files while writing

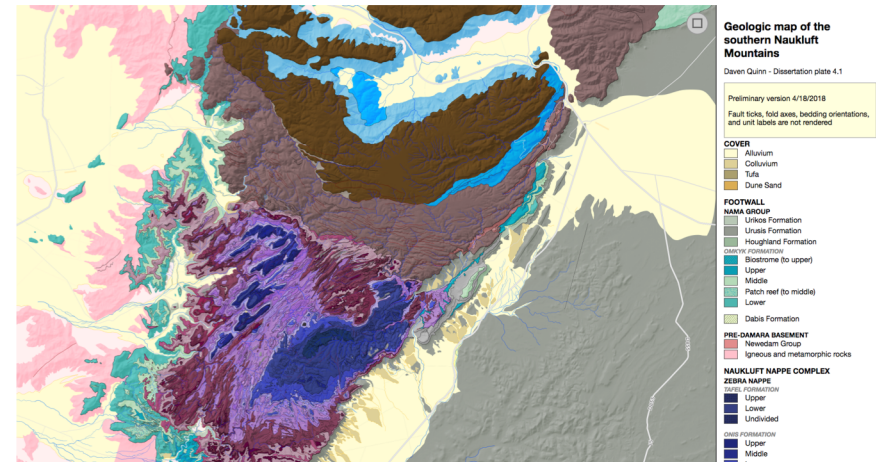


California Institute of Technology  
Research Data Repository

<https://doi.org/10.22002/D1.234>  
<https://doi.org/10.22002/D1.235>  
<https://doi.org/10.22002/D1.236>  
<https://doi.org/10.22002/D1.237>

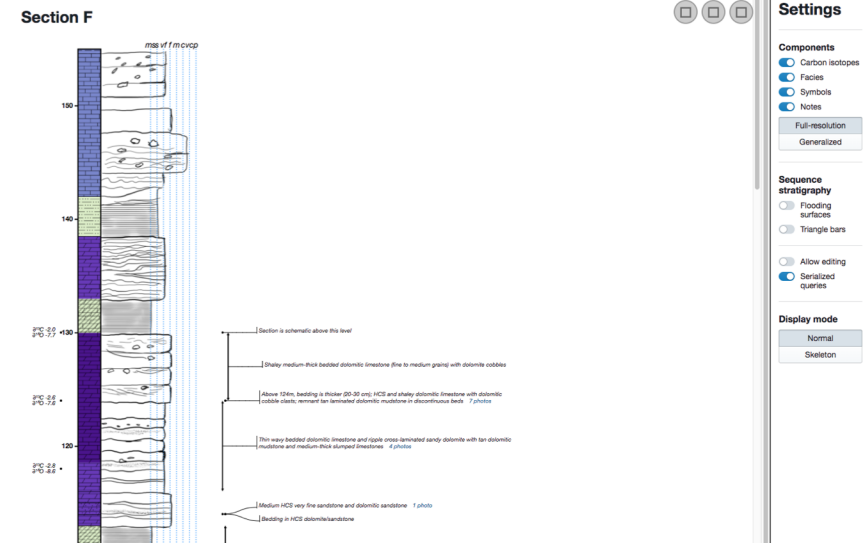
Link in thesis

<https://doi.org/10.7907/Z9NC5Z7H>



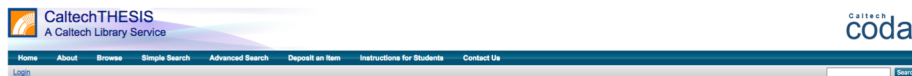
<https://doi.org/10.7907/5exk-mr58>

<https://doi.org/10.22002/D1.946>



<https://doi.org/10.7907/9kva-eq78>

<https://doi.org/10.22002/D1.947>



Engineered Viral Vectors and Developed Tissue Clearing Methods for Single-cell Phenotyping in Whole Organs

Citation  
Chan, Ken Yee (2017) Engineered Viral Vectors and Developed Tissue Clearing Methods for Single-cell Phenotyping in Whole Organs. Dissertation (Ph.D.), California Institute of Technology. doi:10.7907/Z9NC5Z7H. <https://resolver.caltech.edu/10.7907/Z9NC5Z7H>

Abstract  
A central question in biology is how different cell types interact with each other and their native environment to form complex functional systems and networks. Although our ability to investigate this question has considerably expanded from the development of genetically encoded tools, some limitations still persist. For instance, we are limited in our ability to visualize the native three dimensional environments of whole organs. Additionally, it is challenging to efficiently deliver transgene into difficult-to-target areas through direct injections, such as the cardiac ganglia, or broadly distributed networks, such as the myenteric nervous system, which limits our ability to extensively study these areas. Therefore, tools and methods that overcome these limitations are needed. Towards this end, my thesis work has been focused on developing tools for single-cell resolution phenotyping in whole organs. I have been developing tissue clearing technologies to render whole organs transparent for optical interrogation and characterizing viral capsids and engineering viral vectors for noninvasive widespread gene delivery to the central and peripheral nervous system.  
Tissue clearing techniques for three dimensional optical interrogation were invented over a century ago. However, these earlier methods used harsh organic chemicals and failed to retain the tissue's native fluorescence or epitopes. These earlier methods eventually became inapplicable to the hundreds of newly generated transgenic mouse lines that allowed for cell type-specific expression of fluorescent transgenes or to fluorescent labeling techniques, such as immunohistochemistry (IHC). The first part of my dissertation is aimed at addressing these limitations by further developing and standardizing a tissue clearing method that utilizes the available (i) whole organ clearing reagents. This technique, called perfusion assisted agent release in situ (PAAIRS) enables (i) whole organ clearing of soft tissue, (ii) preservation of native fluorescence, and (iii) preservation of epitopes compatible with IHC.

# Use Cases



**bioRxiv**  
beta  
THE PREPRINT SERVER FOR BIOLOGY

HOME | ABO

Search



California Institute of Technology  
Research Data Repository

New Results

## An allosteric theory of transcription factor induction

Manuel Razo-Mejia, Stephanie L. Barnes, Nathan M. Belliveau, Griffin Chure, Tal Einav, Rob Phillips

doi: <https://doi.org/10.1101/111013>

This article is a preprint and has not been peer-reviewed [what does this mean?].

Abstract Info/History Metrics Supplementary material Preview PDF

### Abstract

Allosteric molecules serve as regulators of cellular activity across all domains of life. We present a general theory of allosteric transcriptional regulation that permits quantitative predictions for how physiological responses are tuned to environmental stimuli. To test the model's predictive power, we apply it to the specific case of the ubiquitous simple repression motif in bacteria. We measure the fold-change in gene expression at different inducer concentrations in a collection of strains that span a range of repressor copy numbers and operator binding strengths. After inferring the inducer dissociation constants using data from one of these strains, we show the broad reach of the model by predicting the induction profiles of all other strains. Finally, we derive an expression for the free energy of allosteric transcription factors which enables us to collapse the data from all of our experiments onto a single master curve, capturing the diverse phenomenology of the induction profiles.

<https://doi.org/10.1101/111013>

Paper Website  
on GitHub



The screenshot shows a GitHub repository page. On the left is a navigation menu with links for ABOUT, ANALYSIS, DATA, PEOPLE, and ACKNOWLEDGEMENTS. Below the menu is the Caltech logo and a link to the Phillips Lab GitHub Repo. The main content area features a large green oval containing a handwritten mathematical equation: 
$$\text{Fold-Change} \approx \left(1 + \frac{P}{K} \frac{R}{N_{NS}} e^{-\beta \Delta G_{RA}}\right)^{-1}$$
 Below the equation is the title 'An Allosteric Theory of Transcription Factor Induction' and a description of the website's purpose. At the bottom, there are links for 'Main Text' and 'Supplementary Information'.

Data Files



<https://doi.org/10.22002/D1.224>  
<https://doi.org/10.22002/D1.227>  
<https://doi.org/10.22002/D1.228>  
<https://doi.org/10.22002/D1.229>

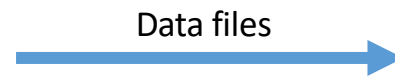
[https://rpgroup-pboc.github.io/mwc\\_induction](https://rpgroup-pboc.github.io/mwc_induction)  
<https://doi.org/10.22002/D1.299>



# Use Case - TCCON



Total Carbon Column Observing Network (TCCON)  
29 Data Collection Sites Around the World



Data Curation and Processing

# Use Case - TCCON



[tcon.ornl.gov](http://tcon.ornl.gov)

TCCON Data Archive HOME GGG2014 GGG2012 GGG2009

Total Carbon Column Observing Network (TCCON)  
The TCCON Data Archive

TCCON is a network of ground-based Fourier Transform Spectrometers recording direct solar spectra in the near-infrared spectral region. From these spectra, accurate and precise column-averaged abundances of CO<sub>2</sub>, CH<sub>4</sub>, N<sub>2</sub>O, HF, CO, H<sub>2</sub>O, and HDO are retrieved. The HF and HDO retrievals are uncalibrated and hence preliminary. Data are updated monthly on the first of the month. The data become publicly available no later than one year after the measurements are recorded, and many sites choose to release their data much sooner.

For the latest TCCON information, please visit the [TCCON Wiki](#). For citation information and our data policy, please see our [Data Use Policy](#). For site-specific information and data analysis descriptions, please read the [Data Description](#). Auxiliary data (column averaging kernels, a priori profiles) are included in the netCDF files provided below. Information on how to use our column averaging kernels and a priori profiles can be found on our [Auxiliary Data](#) page.

A technical report describing the GGG2014 TCCON data version can be found on the [documentation](#) page. Our telluric line list can be downloaded from the [slm](#) page. Our solar line list can be downloaded from the [solar](#) page. A program to generate our a priori profiles can be downloaded from the [a priori](#) page. Please note that the a priori profiles used in the TCCON retrievals are included in the data files below. If you need to produce TCCON a priori profiles for locations and times where there are no TCCON measurements, please use the program linked above.

The TCCON is closely affiliated with the Network for the Detection of Atmospheric Composition Change Infrared Working Group (NDACC-IRWG). In contrast with TCCON, which produces column-averaged dry-air mole fractions, the NDACC produces vertical profiles of the concentrations of many of the same gases and several others. The NDACC website and links to their database can be found at [www.acd.ucar.edu/irwg](http://www.acd.ucar.edu/irwg).

Sign up to the TCCON Users email list to get email updates on TCCON data releases.  
Note that the website is self-signed; you can safely add an exception.

[Login for TCCON Partners](#)

## Private data files

- Sites
- Ascension Island
- [0ae20120522\\_20120831.nc](#)
  - [0ae20130317\\_20130618.nc](#)
  - [0ae20130911\\_20131229.nc](#)
  - [0ae20140108\\_20140716.nc](#)
  - [0ae20140717\\_20141019.nc](#)
  - [0ae20141021\\_20141231.nc](#)
  - [0ae20150101\\_20150310.nc](#)
  - [0ae20150311\\_20150409.nc](#)
  - [0ae20150410\\_20150630.nc](#)
  - [0ae20150701\\_20150926.nc](#)
  - [0ae20151005\\_20151218.nc](#)

## Public data files

### Index of /2014Public/ascension01

Name	Size	Date Modified
<a href="#">[parent directory]</a>		
<a href="#">README.txt</a>	11.8 kB	10/20/14, 5:00:00 PM
<a href="#">ae20120522_20161221.public.nc</a>	10.1 MB	5/31/17, 5:25:00 PM

Automatically released 1x/month

## Departmental Server at Caltech

Login for TCCON Partners TCCON Data Archive GGG2014 GGG2012 GGG2009

### Total Carbon Column Observing Network (TCCON)



TCCON is a network of ground-based Fourier Transform Spectrometers recording direct solar spectra in the near-infrared spectral region. From these spectra, accurate and precise column-averaged abundances of CO<sub>2</sub>, CH<sub>4</sub>, N<sub>2</sub>O, HF, CO, H<sub>2</sub>O, and HDO are retrieved and reported here. A technical report describing the retrievals is found [here](#); solar and telluric spectral line lists used in the retrievals are publicly available.

Data in netCDF format are publicly available no later than one year after the spectra are recorded; many sites release their data earlier. Citation and data use requirements are included in the license associated with each record. Column averaging kernels and a priori profiles are included in the files. Information on how to use these can be found [here](#). To produce TCCON a priori profiles for locations and times where there are no TCCON measurements, a stand-alone program can be [downloaded](#).

Sign up to the TCCON Users email list to get email updates on TCCON data releases.

[tcondata.org](http://tcondata.org)

CaltechDATA

## Migration

TCCON data from Park Falls (US), Release GGG2014.R1

Dataset

2017-09-27

CaltechDATA

Download Edit

**Details**

**Authors**  
Wernberg, P. O. California Institute of Technology, Pasadena, CA (US) 0000-0002-6126-3854 ORCID  
Roehl, C. M. California Institute of Technology, Pasadena, CA (US) 0000-0001-5383-8462 ORCID  
Wunch, D. California Institute of Technology, Pasadena, CA (US) 0000-0002-4024-0377 ORCID  
Toon, G. C. Jet Propulsion Laboratory, Pasadena, CA (US)  
Blavier, J.-F. Jet Propulsion Laboratory, Pasadena, CA (US) 0000-0001-1808-8316 ORCID  
Washenfelder, B. University of Colorado, NOAA, Boulder, CO (US) 0000-0002-8106-3702 ORCID  
Koppel-Aleks, G. University of Michigan, Ann Arbor, MI (US) 0000-0003-2119-0044 ORCID  
Allen, N. T. Harvard University, Cambridge, MA (US) 0000-0002-7528-8605 ORCID  
Ayers, J. Wisconsin Educational Communications Board, Park Falls, WI (US)

**Contributors**  
Hosting Institution California Institute of Technology, Pasadena, CA (US)  
Data Curator Roehl, C. M. California Institute of Technology, Pasadena, CA (US) 0000-0001-5383-8462 ORCID  
Contact Person Paul Wernberg [wernberg@jpl.caltech.edu](mailto:wernberg@jpl.caltech.edu)

**Description**  
Abstract:  
The Total Carbon Column Observing Network (TCCON) is a network of ground-based Fourier Transform Spectrometers that record direct solar absorption spectra of the atmosphere in the near-infrared. From these spectra, accurate and precise column-averaged abundances of atmospheric constituents including CO<sub>2</sub>, CH<sub>4</sub>, N<sub>2</sub>O, HF, CO, H<sub>2</sub>O, and HDO, are retrieved. This data set contains observations from the TCCON station at Park Falls, U.S.A.

**Publication Date**  
2017-09-27

**Subjects**  
atmospheric trace gases, CO2, CH4, CO, N2O, column-averaged dry-air mole fractions, remote sensing, FTIR spectroscopy, TCCON

**DOI**  
10.14291/tcon-ggg2014-parkfalls01.R1

**Version**  
GGG2014.R1

**Format**  
application/x-netcdf

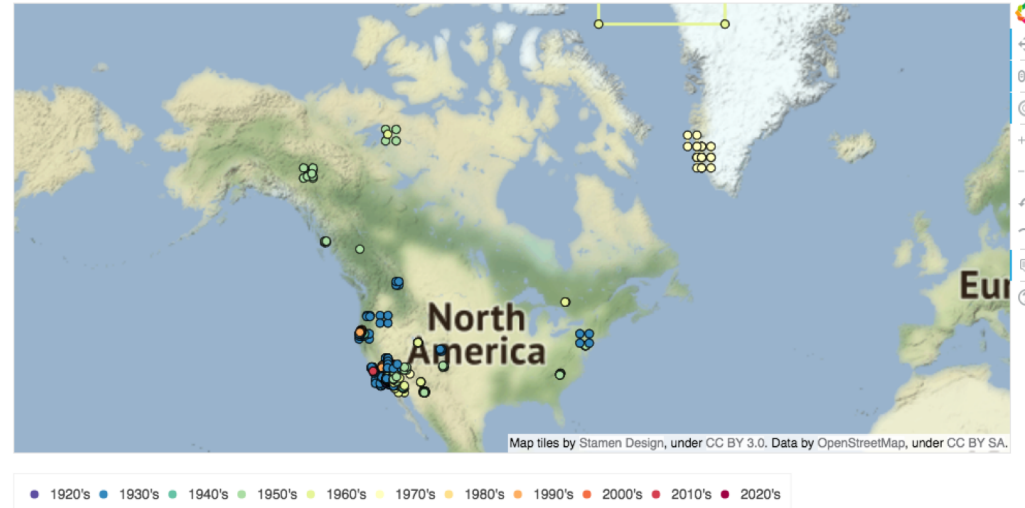
<https://doi.org/10.14291/tcon.ggg2014.parkfalls01.R1>



### Caltech Division of Geological and Planetary Sciences Theses

This map shows the coordinates of content in CaltechDATA associated with theses from the Geological and Planetary Science Division at Caltech. Data included from historic theses are supplemental pocket contents such as maps and drawings.

Scrolling inside the map will zoom and dragging will move the map. Click on any point or bounding box to see the original item in CaltechDATA.



#### Want your thesis to show up on the map?

Upload files associated with your thesis to CaltechDATA and include a geolocation point or area. You'll also have to include the keywords 'gps' and 'thesis' in the record. If you run into any problems just send us an email.

#### Did you complete your thesis in the Caltech GPS Division?

We haven't been able to assign locations for every thesis. Send us an email and we can get your thesis on the map.

#### Want to improve this map?

The code to generate the map is available on GitHub and we accept pull requests for improvements.

<http://maps.library.caltech.edu/>

<https://doi.org/10.22002/D1.856>

# Caltech Library Data Management Services

- Want to chat about data issues?
- Data management plan development
- Consultations on storage technologies or file organization

[data@caltech.edu](mailto:data@caltech.edu)

[tmorrell@caltech.edu](mailto:tmorrell@caltech.edu)

626-395-3827