# The Geoscience Paper of the Future:
## Practical Guidelines for Adopting Digital Scholarship, Reproducible Research, and Open Science

**Yolanda Gil**

Information Sciences Institute and
Department of Computer Science
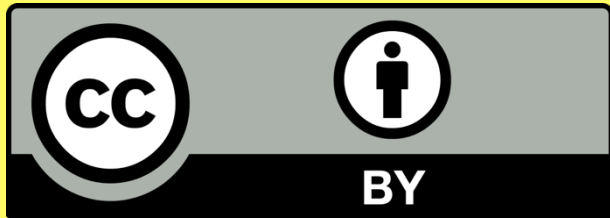University of Southern California

gil@isi.edu

http://www.ontosoft.org/gpf

ICER-1440323
ICER-1343800

**InGeO** Integrated Geoscience Observatory

ICER-1541057

EarthCube!

# Acknowledgments

ICER-1440323
ICER-1343800

InGeO
Integrated
Geoscience
Observatory

ICER-1541057

# Problems with Current Practice

- ★ Data is often not made available in publications
  - ★ Lack of reproducibility

*Nature Genetics* **41**, 149 - 155 (2009)
Published online: 28 January 2008 | doi:10.1038/ng.295

## Repeatability of published microarray gene expression analyses

scientists. Here we evaluated the replication of data analyses in 18 articles on microarray-based gene expression profiling published in *Nature Genetics* in 2005–2006. One table or figure from each article was independently evaluated by two teams of analysts. We reproduced two analyses in principle and six partially or with some discrepancies; ten could not be reproduced. The main reason for failure to reproduce was data unavailability, and discrepancies were mostly due to incomplete data annotation or specification of data processing and analysis.

- ★ Data made available through URLs that are not persistent
  - ★ URL does not resolve (i.e., ''rotten'')

PLOS ONE | DOI:10.1371/journal.pone.0115253    December 26, 2014

RESEARCH ARTICLE

## Scholarly Context Not Found: One in Five Articles Suffers from Reference Rot

Martin Klein[1]*, Herbert Van de Sompel[1], Robert Sanderson[1], Harihar Shankar[1], Lyudmila Balakireva[1], Ke Zhou[2], Richard Tobin[2]

We analyze a vast collection of articles from three corpora that span publication years 1997 to 2012. For over one million references to web resources extracted from over 3.5 million articles, we observe that the fraction of articles containing references to web resources is growing steadily over time. We find one out of five STM articles suffering from reference rot, meaning it is impossible to revisit the web context that surrounds them some time after their publication. When only considering STM articles that contain references to web resources, this fraction increases to seven out of ten.

# Publishers Are Changing: Guidelines for Authors



**nature research**

Data availability statements and data citations policy: guidance for authors

Policy summary

All manuscripts reporting original research must include a data availability statement. Authors are also encouraged to include formal citations to datasets in article reference lists where deposited datasets are assigned Digital Object Identifiers (DOIs) by a data repository.

nature.com > scientific data

SCIENTIFIC DATA

nature.com

protocol exchange



**PLOS | ONE**

**Availability of Software**

PLOS supports the development of open source software and believes that, for submissions appropriate open source standards will ensure that the submission conforms to (1) our requirements another researcher can reproduce the experiments described, (2) our aim to promote openness PLOS journals can be built upon by future researchers. Therefore, if new software or a new that the software conforms to the Open Source Definition, have deposited the following three submission as Supporting Information:

- The associated source code of the software described by the paper. This should be licensed under a suitable license such as BSD, LGPL, or MIT (see http://www.ope commercial software such as Mathematica and MATLAB does not preclude a paper preferred.
- Documentation for running and installing the software. For end-user applications prerequisite; for software libraries, instructions for using the application program inter
- A test dataset with associated control parameter settings. Where feasible, result test data should not have any dependencies — for example, a database dump.

Acceptable archives should provide a public repository of the described software. The code for creating user accounts, logging in or otherwise registering personal details. The repositor more than 1,000 projects. Examples of such archives are: SourceForge, Bioinformatics.Org, Savannah, GitHub and the Codehaus. Authors should provide a direct link to the deposited s



**COPDESS**

Coalition on Publishing Data in the Earth and Space Sciences

**COPDESS Suggested Author Instructions and Best Practices for Journals**

The Coalition on Publishing Data in the Earth and Space Sciences (COPDESS) develops and recommends best practices for journal author instructions around data and identifiers as a resource to the community. These best practices are consistent with and based on the COPDESS Statement of Commitment and have been developed with guidance from participants in COPDESS.
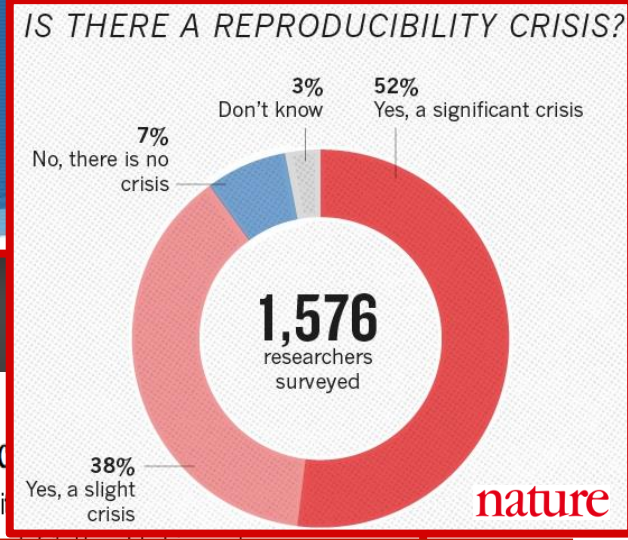
Data Policy Statement
Data Citation
Sample Citation and Identification
Crossref Funder Registry
ORCIDs
Presentations on Best Practices

# Reproducibility

**IS THERE A REPRODUCIBILITY CRISIS?**

3% Don't know
52% Yes, a significant crisis
7% No, there is no crisis
1,576 researchers surveyed
38% Yes, a slight crisis

nature

**Illuminating the black box**

Note to biologists: submissions to *Nature* should contain complete descriptions of materials and reagents used.

Reporting Checklist For Life Sciences Articles

This checklist is used to ensure good reporting standards and please read Reporting Life Sciences Research

nature

**Methodology**

Friday, December 2, 2011 As of 12:00 AM   New York   43° | 34°

THE WALL STREET JOURNAL. | HEALTH

HEALTH INDUSTRY | DECEMBER 2, 2011

Scientists' Elusive Goal: Reprodu

n September, Bayer published a study describing how i

COMPUTER SCIENCE

**Accessible Reprodu**

Jill P. M

Science

POLICYFORUM

putation in research grows, needed to expand recording,

**A Biostatistic Paper Alleges Potential Harm To Patients In Two Duke Clinical Studies**

statistics journals   r sensational s. The most recent issue of the Annals of Applied Statistics is an

**Human lives**

lleges that cancer patients

**Reliability**

The New York Times

NYTimes: Home - Site Index - Archive - Help

**Nobel Laureate Retracts Two Paper**

B. KENNETH CHANG

The New York Times

Retracted Scientific Studies: A Growing List

RETRACTED

**Scientific integrity**

新语丝

New Threads

**No Cure**

When Bayer tried to replicate results of 67 studies published in academic journals, nearly two-thirds failed.

Fully replicated 20.9%
Partially replicated 11.9%
Not replicated 64.2%
Not applicable 3.0%

Source: Nature Reviews Drug Discove

**Financial**

**Trust**

GLOBAL WARMING SCIENCE FICTIO

# Government Agencies Are Changing: Scientific Integrity and Open Science

Office of Science and Technology Policy

## Scientific Integrity

On December 17, 2010, OSTP Director John P. Holdren issued a Memorandum for the Heads of Executive Departments and Agencies on the subject of Scientific Integrity.

- Read the blog
- Read the December 17, 2010 Memorandum (pdf)
- Read the President's March 9, 2009 Memorandum
- Read sample communications policy language (pdf)

https://obamawhitehouse.archives.gov/administration/eop/ostp/library/scientificintegrity

- Department of Agriculture (pdf)
- Department of Commerce (pdf)
  - National Institute of Standards and Technology (pdf)
  - National Oceanic and Atmospheric Administration
- Department of Defense (pdf)
- Department of Education (pdf)
- Department of Energy (pdf)
- Department of Health and Human Services (pdf)
  - Centers for Disease Control and Prevention (pdf)
  - Food and Drug Administration
  - National Institutes of Health (pdf)
- Department of Homeland Security (pdf)
- Department of the Interior (pdf)
- Department of Justice (pdf)
- Department of Labor (pdf)
- Department of State (pdf)
- Department of Transportation
- Department of Veteran Affairs (pdf)
- United States Agency for International Development
- Environmental Protection Agency (pdf)
- Marine Mammal Commission (pdf)
- National Aeronautics and Space Administration (pdf)
- National Science Foundation (pdf)
- Office of the Director of National Intelligence (pdf)

# Growing Importance of Scientific Integrity and Reproducibility

# Science is Changing: Sharing, Open, Credit

# Universities are Changing: Major Initiatives in Data Science

# Core Recommendations for Scientific Publications

**Reproducible Research**

**Open Science**

**Digital Scholarship**

# 1) Reproducible Research



**Datasets**

**Software**

**Workflow**

**Experimental Design**

**Open Science**

**Shared repositories**

**Persistent unique identifiers**

**Licenses**

# 3) Digital Scholarship

Digital Scholarship

**Citation**

**Metadata**

Garijo, Daniel; Xie, Lei; Zhang, Yinliang; Gil, Yolanda; Xie, Li; Kinnings, Sarah; Bourne, Phil (2013) Highly connected drug file figshare.
http://dx.doi.org/10.6084/m9.figshare.776887
Retrieved 11:05, Feb 20, 2015 (GMT)

Authors

Date of publication

Time of retrieval

Persistent unique identifier

Name

Repository

Garijo, Daniel;Xie, Lei; Zhang, Yinliang; Gil, Yolanda; Xie, Li (2013) Tool for computing anomalies, GitHub. V.1
http://dx.doi.org/10.5281/zenodo.18765
Retrieved 11:05, Feb, 15, 2015 (GMT)

BIBLIOMETRICS AND CITATION ANALYSIS
From the *Science Citation Index* to Cybermetrics

Journal of open research software

Version

# Geoscience Paper of the Future

## Modern Paper

**Text:**
Narrative of the method, some data is in tables, figures/plots, and the software used is mentioned

**Data:**
Include data as supplementary materials and pointers to data repositories

## Open Science

**Sharing:**
Deposit data and software (and provenance/workflow) in publicly shared repositories

**Open licenses:**
Open source licenses for data and software (and provenance/workflow)

**Metadata:**
Structured descriptions of the characteristics of data and software (and provenance/workflow)

## Reproducible Publication

**Software:**
For data preparation, data analysis, and visualization

**Provenance and methods:**
Workflow/scripts specifying dataflow, codes, configuration files, parameter settings, and runtime dependencies

## Digital Scholarship

**Persistent identifiers:**
For data, software, and authors (and provenance/workflow)

**Citations:**
Citations for data and software (and provenance/workflow)

# The Geoscience Papers of the Future (GPF) Initiative

1. A Special Issue of a journal in all geoscience areas that includes only geoscience papers of the future



**Special Section: Geoscience Papers of the Future**

1. Training sessions for geoscientists to learn best practices in software and data sharing, provenance documentation, and scholarly publication

# GPF Pioneer Authors

**Cedric David**, NASA/JPL
Hydrology modeling

**Ibrahim Demir**, U. of Iowa
Hydrology sensor networks

**R. W. Fulweiler**, Boston U.
Biogeochemistry in marine ecology

**J. Goodall/B. Essawy**, U.
Virginia, Hydrology/visualization

**Leif Karlstrom**, U. Oregon
Volcanic vent clustering

**Kyo Lee**, NASA/JPL
Regional climate modeling

**Heith Mills**, U. Houston
Geochemistry, marine biology

**Ji-Hyun Oh**, USC
Tropical meteorology

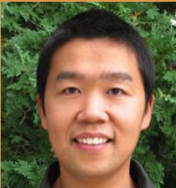**Suzanne Pierce**, UT Austin
Hydrogeology for decision support

**Allen Pope**, U. Colorado
Glaciology

**Mimi Tzeng**, Dauphin Island
Sea Lab, Ocean fisheries

**Sandra Villamizar**, UC Merced
River ecohydrology

**Xuan Yu**, U. Delaware
Hydrologic modeling

# Published Articles

## Geoscience Paper of the Future

### Modern Paper

**Text:**
Narrative of the method, some data is in tables, figures/plots, and the software used is mentioned

**Data:**
Include data as supplementary materials and pointers to data repositories

### Open Science

**Sharing:**
Deposit data and software (and provenance/workflow) in publicly shared repositories

**Open licenses:**
Open source licenses for data and software (and provenance/workflow)

**Metadata:**
Structured descriptions of the characteristics of data and software (and provenance/workflow)

### Reproducible Publication

**Software:**
For data preparation, data analysis, and visualization

**Provenance and methods:**
Workflow/scripts specifying dataflow, codes, configuration files, parameter settings, and runtime dependencies

### Digital Scholarship

**Persistent identifiers:**
For data, software, and authors (and provenance/workflow)

**Citations:**
Citations for data and software (and provenance/workflow)

---

*"Towards the Geoscience Paper of the Future: Best Practices for Documenting and Sharing Research from Data to Software to Provenance"* Gil et al, Earth and Space Science, 2016.
http://dx.doi.org/10.1002/2015EA000136

---

- [David et al 2015]: 10 years of hydrology model software
- [Yu et al 2015]: Model coupling for surface/subsurface flow
- [Essawy et al 2015]: Hydrology workflows for reproducibility
- [Pope et al 2015]: Estimate subglaciar lake depth from imagery
- [Fulweiler et al 2016]: Long-term estuary data & products
- [Tzeng et al 2016]: Data processing for ocean observatory
- [Demir et al 2017]: Sensor network for flood monitoring
- *<more in process>*

# An Example

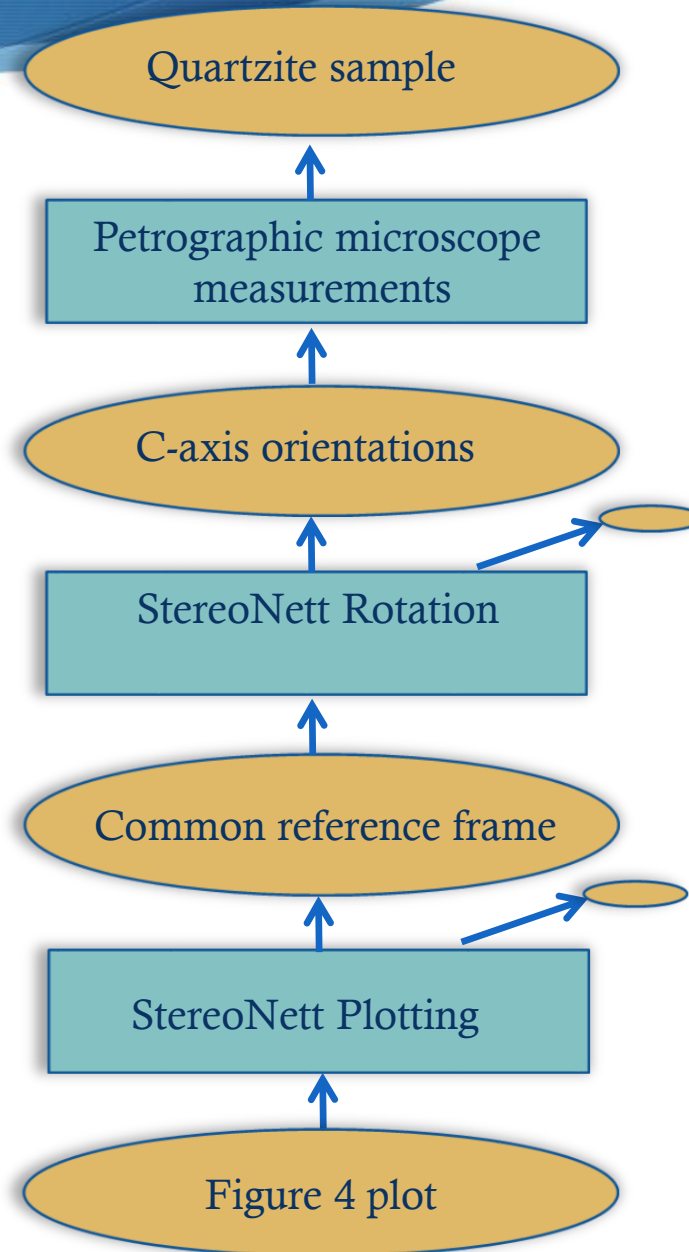**Understanding kinematic data from the Hellerman thrust zone**

**Jade Silverstein**

**[…] We took a quartzite sample from the Hellerman thrust zone, and cut 3 thin sections. We measured c-axis orientations using a petrographic microscope. We rotated to a common reference frame using Duyster's StereoNett program. We plotted the data on lower hemisphere, equal area projections using Duyster's StereoNett program, shown in Figure 4. […]**

# An Example

**Understanding kinematic data from the Hellerman thrust zone**

**Jade Silverstein**

**[…] We took a quartzite sample from the Hellerman thrust zone, and cut 3 thin sections. We measured c-axis orientations using a petrographic microscope. We rotated to a common reference frame using Duyster's StereoNett program. We plotted the data on lower hemisphere, equal area projections using Duyster's StereoNett program, shown in Figure 4. […]**
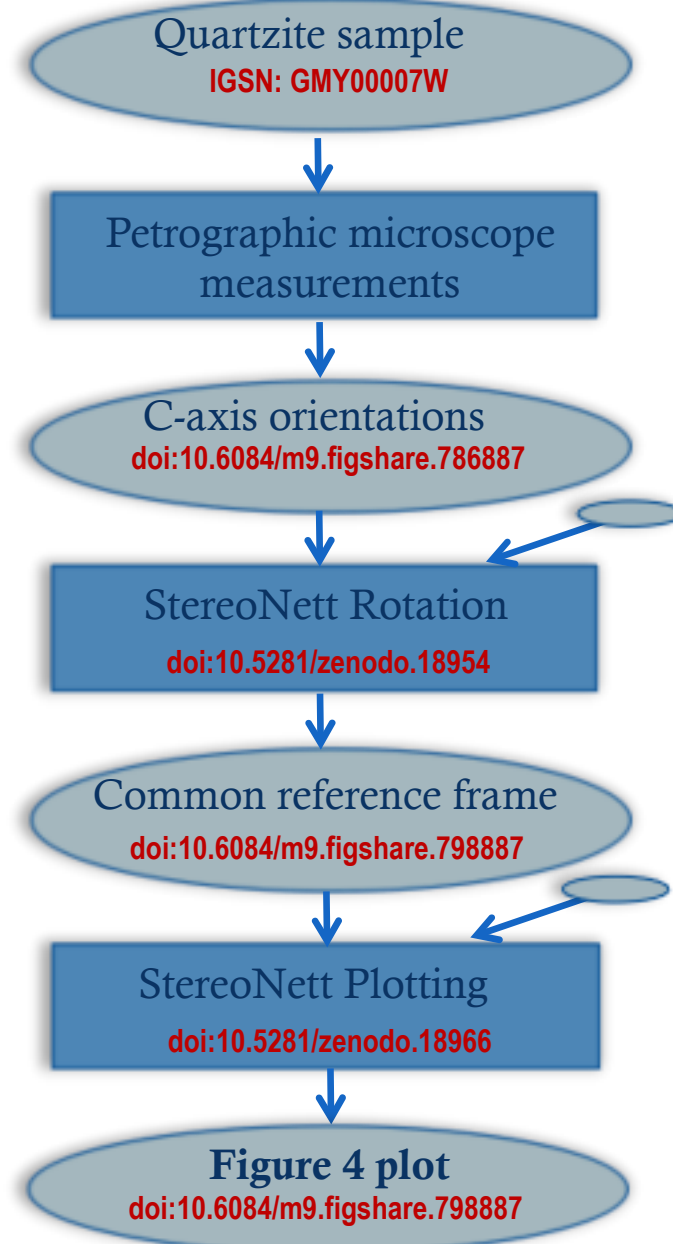
Quartzite sample

Petrographic microscope measurements

C-axis orientations

StereoNett Rotation

Common reference frame

StereoNett Plotting

Figure 4 plot

# An Example

**Understanding kinematic data from the Heller thrust zone (doi:10.1016/j.ess.2009.08.012)**

**Jade Silverstein (orcid.org/0000-0001-8455-8431)**

**[…] We took a quartzite sample (IGSN: GMY00007W) from the Heller thrust zone, and cut 3 thin sections. We measured c-axis orientations (doi:10.6084/m9.figshare.786887) using a petrographic microscope. We rotated to a common reference frame (doi:10.6084/m9.figshare.798887) using Duyster's StereoNett program (doi:10.5281/zenodo.18954). We plotted the data on lower hemisphere, equal area projections (doi:10.6084/m9.figshare.798887) using Duyster's StereoNett program (doi:10.5281/zenodo.18966), shown in Figure 4. The provenance is shown in Fig 5. […]**

Quartzite sample
IGSN: GMY00007W

Petrographic microscope measurements

C-axis orientations
doi:10.6084/m9.figshare.786887

StereoNett Rotation
doi:10.5281/zenodo.18954

Common reference frame
doi:10.6084/m9.figshare.798887

StereoNett Plotting
doi:10.5281/zenodo.18966

Figure 4 plot
doi:10.6084/m9.figshare.798887

# Modern Scientific Articles

**Traditional Published Articles**

**Text**:
Narrative of method,
the data is in tables, figures/plots,
the software used is mentioned

**Modern Published Articles**

**Text:**
Narrative of method,
the data is in tables, figures/plots,
the software used is mentioned

**Data:**
Supplementary materials,
pointers to data repositories

**NOT published,**
**loosely recorded:**

**Software:**
scripted codes + manual steps +
documentation in notes/emails

# Reproducible Articles

## Modern Published Articles

**Text:**
Narrative of method,
the data is in tables, figures/plots,
the software used is mentioned

**Data:**
Supplementary materials,
pointers to data repositories

**NOT published,
loosely recorded:**

**Software:**
scripted codes + manual steps +
documentation in notes/emails

## Reproducible Publications

**Text:**
Narrative of method,
the data is in tables, figures/plots,
the software used is mentioned

**Data:**
Supplementary materials,
pointers to data repositories

**Software:**
Data preparation,
data analysis, and visualization

**Provenance and Workflow:**
Workflow/scripts describing
dataflow, codes, and parameters
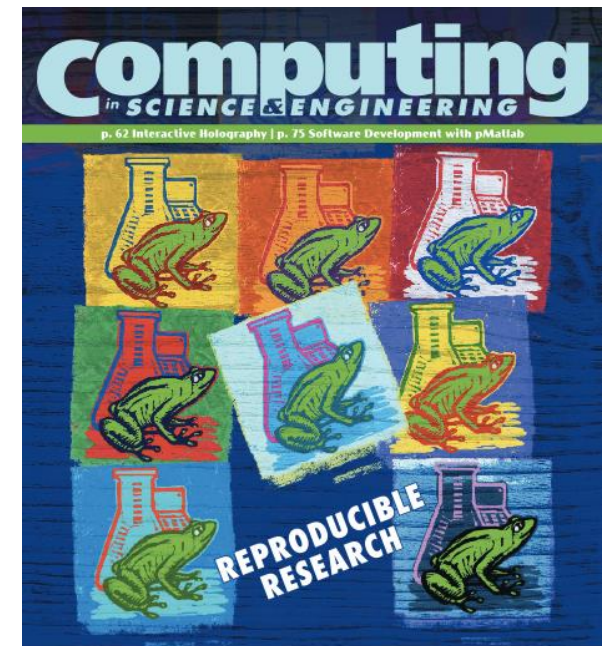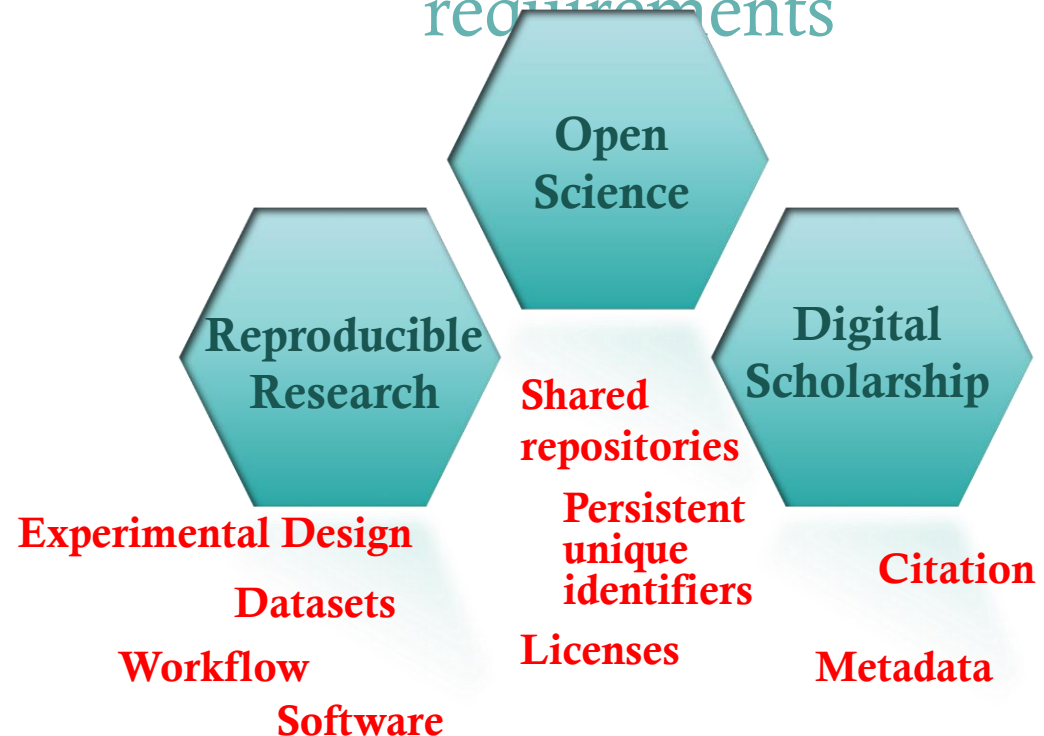
# Reproducible Publications and Executable Papers

MADAGASCAR

$$Sweave = R \cdot \LaTeX$$

IP[y]: Notebook

Science

Data Replication and Reproducibility

AAAS

computing
in SCIENCE & ENGINEERING

p. 62 Interactive Holography | p. 75 Software Development with pMatlab

REPRODUCIBLE
RESEARCH

IEEE

January/February 2009, Vol. 11, No. 1

# Beyond Reproducible Publications

## Reproducible Publications

**Text:**
Narrative of method,
the data is in tables, figures/plots,
the software used is mentioned

**Data:**
Supplementary materials,
pointers to data repositories

**Software:**
Data preparation,
data analysis, and visualization

**Provenance and methods:**
Workflow/scripts describing
dataflow, codes, and parameters

The Geoscience Paper of the Future has further requirements

**Open Science**

**Reproducible Research**

**Digital Scholarship**

Shared repositories

Persistent unique identifiers

Experimental Design

Datasets

Workflow

Software

Licenses

Citation

Metadata

Research paper

# Simulating electron and ion temperature in a global ionosphere thermosphere model: Validation and modeling an idealized substorm

Jie Zhu [a,*], Aaron J. Ridley [a], Yue Deng [b]

[a] Department of Atmospheric, Oceanic and Space Sciences, University of Michigan, Ann Arbor, MI, USA
[b] Department of Physics, University of Texas, Arlington, USA

A R T I C L E   I N F O

A B S T R A C T

Electron and ion temperatures control many chemical and physical processes in th
mosphere system. Recently, improved electron and ion energy equations were i
Global Ionosphere Thermosphere Model (GITM). The source energy of the electro
includes thermal conduction, heating due to photoionization, elastic collisions wi
inelastic collisions with neutrals, auroral precipitation, and heat flux from inner r
source terms in the ion temperature ($T_i$) equation include thermal conduction, and e
electrons and neutrals. The new implementation of $T_e$ improved the ionospheric de
high latitudes with respect to IRI. The improved GITM also reproduced the diurnal v
observed by incoherent scatter radars at low and middle latitudes. The model was us
idealized substorm statistically described by Clausen et al. (2014). It was found that

# Computational Aspects of the Paper

- ★ Comparison of a new version of the Global Ionosphere Thermosphere Model (GITM) with the previous version

- ★ Comparisons with the IRI model (Bilitza et al 2014) for several sites

- ★ The model was used to investigate an idealized substorm

# Data

## From Paper

- Fig. 1 shows comparisons … using the old model (left), the new model (middle) and IRI (right) at an altitude of 400 km at 00:00 UT on December 23, 2012.
  - Missing (though provided by CCMC):
    - Brightness of sun: NOAA
    - Strength of aurora: NOAA
    - Electric fields: NASA
- Investigate the ionospheric response to an idealized substorm, which was the same as Substorm 4 investigated by Liu and Ridley (2015). The prototypical substorm was constructed based on the superposed epoch variations of IMF Bz and HP during substorms using 5-years of Challenging Minisatellite Payload (CHAMP) (Reigber et al., 2002) satellite data (Clausen et al., 2014).
  - The (Liu and Ridley 2015) paper has the substorm data and plots for it

## Best Practices

- **All input data should be in a public repository (as well as any important <u>intermediate data</u>)**
  - **Community repositories (Madrigal), university (Dataverse), other (zenodo)…**
- **A unique identifier (DOI) should be assigned to each dataset**
- **Basic metadata should be attached**
- **A license should be specified for each dataset**
  - **Creative Commons**
    - **Recommendation: CC-BY**
- **Data should be cited in-line, the citation should be included in the references**

# Software

## From Paper

* Tables 1 and 2 present a comparison of the implementation of Te and Ti between the old and new model
    * Missing:
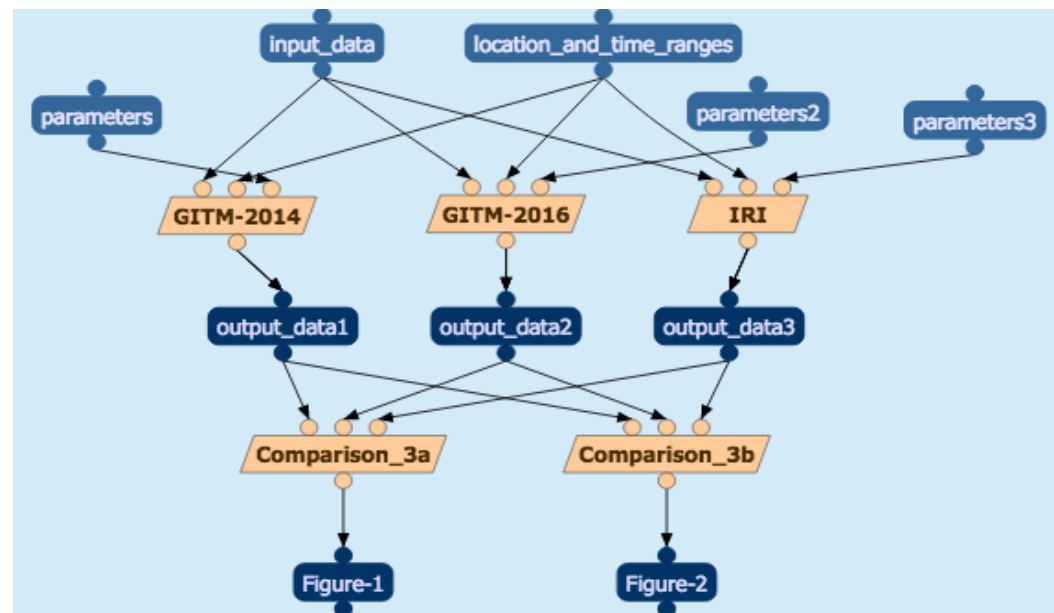        * Pointers to new model and previous model software versions

* Fig. 1 shows comparisons of Ne (top), Te (middle) and Ti (bottom) using the old model (left), the new model (middle) and the International Reference Ionosphere (IRI) model (Bilitza et al., 2014) (right)
    * Missing:
        * Pointer to ISI software version

## Best Practices

* **All software should be in a public repository**
    * **Community repository**
        * **GitHub, or zenodo**
* **Unique identifier (DOI) should be assigned to each software version**
* **Basic metadata should be attached**
    * **See www.ontosoft.org**
* **License should be specified**
    * **See www.opensource.org**
        * **Recommend: Apache v2, MIT**
* **Software should be cited in-line (including version), the citation should be included in the references**

COMMUNITY
COORDINATED
MODELING
CENTER

Related Links | Frequently Asked Questions | Community Feedback | Downloads | Sitemap

NASA    NSF

| About | Models at CCMC | Request A Run | View Results | Instant Run | Metrics and Validation | Education | R2O Support |

# International Reference Ionosphere - IRI-2007

**This page enables the computation and plotting of IRI parameters: electron and ion (O+, H+, He+, O2+, NO+) densities, total electron content, electron, ion and neutral (CIRA-86) temperatures, equatorial vertical ion drift and others.**

Go to the IRI description

Help

**\* Select Date and Time**
**Year**(1960-2017): 2000

**Note:**If date is outside the Ap index range (1960/02/14-2017/03/17 ),then STORM model will be turned off.
**Month:** January ⌄ **Day(1-31):** 01
**Time** Universal ⌄ **Hour of day** (e.g. 1.5): 1.5
**\* Select Coordinates**
**Coordinates Type** Geographic ⌄
**Latitude**(deg.,from -90. to 90.): 50.    **Longitude**(deg.,from 0. to 360.) 40.
**Height** (km, from 60. to 2000.): 100.
**\* Select a Profile type and its parameters:**
Height,km [ 60. - 2000.] ⌄ **Start** 100.    **Stop** 2000.    **Stepsize** 50.

# Global Ionosphere Thermosphere Model (GITM)

**CCMC Services available for GITM**
Request a Run
View Request Results

**Model Developer(s)**
A.J. Ridley et al.
Department of Atmosphere, Oceanic and Space Sciences, University of Michigan

**Model Description**
GITM is a 3-dimensional spherical code that models the Earth's thermosphere and ionosphere system using a stretched grid in latitude and altitude. The number of grid points in each direction can be specified, so

12/21/2002 Time = 07:00:03 UT z= 400. km

50.

★ Accessing older versions?

★ Identifying older versions uniquely?

# Workflow

## From Paper

★ Fig. 1 shows comparisons of Ne (top), Te (middle) and Ti (bottom) using the old model (left), the new model (middle) and the International Reference Ionosphere (IRI) model (Bilitza et al., 2014) (right)



## Best Practices

★ **Sketch a workflow diagram**
  ★ **Software invocations**
  ★ **Dataflow connections**
★ **Include workflow diagram in figure or supplementary materials**

# Provenance of Results

## From Paper

Fig. 5 shows a comparison between GITM, IRI and measurement by the Arecibo ISR from 100 km to 650 km on April 13th and 14th, 2013. The missing data in the observations were filled by an altitudinal linear interpolation.

**ARECIBO**



**JICAMARCA**



**SONDRESTROM**



## Best Practices

★ **Specification of all input data and parameters**

★ **Publication of intermediate and final results facilitates reproducibility**

# Author Checklist

**1** Data accessibility

**2** Data documentation

**3** Software accessibility

**4** Software documentation

**5** Methods documentation

**6** Provenance documentation

**7** Author identification

★ **For datasets**, the paper should include one or more citations, specifying the authors, the site where they are described and can be accessed, the repository, and the license.

★ **For software**, the paper should include one or more citations, specifying the authors, the site where it is described and can be accessed, the repository, and the license.

★ **For provenance and workflow**, the paper should include figures and traces, and if available the citations mentioning the authors, site to access them, the repository, and the license.

★ **For authors**, each should have a unique identifier (e.g., ORCID)

# Incorporate GPF Best Practices Into Your Work



- Easier to track research products, train new lab members, build on prior work, etc.
- Making a paper into a GPF is then very straightforward

# Why Learn to Write a Scientific Paper of the Future

1. Practice **open science and reproducible research**
2. **Get credit** for all your research products
   - ★ Citations for software, data, samples, …
3. **Increase citations** of your papers
4. Write impressive **Data Management Plans**
5. **Extend your CV** with data and software sections
6. Improve the **management of your research assets**
7. **Reproduce** your work from years ago and build on it
8. Address new **funder and journal requirements**
9. Attract **transformative students**
10. Demonstrate **leadership** by stepping into the future

# Recommendations from Scientific Societies

**CRA** Computing Research Association

**Incentivizing Quality and Impact: Evaluating Scholarship in Hiring, Tenure, and Promotion**

*"The field benefits when researchers build on each other's work. To do so, requires that research advances be accompanied by discussion of methods, comparisons with related work, inclusion of supporting data and proofs, access to artifacts, and other details. Certain publication formats and review processes, however, encourage practices inconsistent with these elements of good scholarship. **Length restrictions often are satisfied by omitting critical content, which hinders reproducing the results, understanding their novelty, or delimiting a contribution's applicability. The omission of supporting data and proofs, also common practice, hobbles efforts to validate or extend the work**."*

# For More Information

## Geoscience Paper of the Future

### Modern Paper

**Text:**
Narrative of the method, some data is in tables, figures/plots, and the software used is mentioned

**Data:**
Include data as supplementary materials and pointers to data repositories

### Reproducible Publication

**Software:**
For data preparation, data analysis, and visualization

**Provenance and methods:**
Workflow/scripts specifying dataflow, codes, configuration files, parameter settings, and runtime dependencies

### Open Science

**Sharing:**
Deposit data and software (and provenance/workflow) in publicly shared repositories

**Open licenses:**
Open source licenses for data and software (and provenance/workflow)

**Metadata:**
Structured descriptions of the characteristics of data and software (and provenance/workflow)

### Digital Scholarship

**Persistent identifiers:**
For data, software, and authors (and provenance/workflow)

**Citations:**
Citations for data and software (and provenance/workflow)

**GPF recommended best practices:**
http://dx.doi.org/10.1002/2015EA000136

**Special issue:**
http://tinyurl.com/ess-gpf

**Training materials:**
http://dx.doi.org/10.5281/zenodo.15920

GEOSCIENCE PAPERS OF THE FUTURE
OntoSoft

OntoSoft

InGeO
Integrated Geoscience Observatory

NSF

EarthCube!

ICER-1440323
ICER-1343800

ICER-1541057

# EXTRA SLIDES:
## OVERVIEW OF GPF RECOMMENDATIONS AND AUTHOR CHECKLIST

# Author Checklist

1. Data accessibility
2. Data documentation
3. Software accessibility
4. Software documentation
5. Methods documentation
6. Provenance documentation
7. Author identification

★ **For datasets**, the paper should include one or more citations, specifying the authors, the site where they are described and can be accessed, the repository, and the license.

★ **For software**, the paper should include one or more citations, specifying the authors, the site where it is described and can be accessed, the repository, and the license.

★ **For provenance and workflow**, the paper should include figures and traces, and if available the citations mentioning the authors, site to access them, the repository, and the license.

★ **For authors**, each should have a unique identifier (e.g., ORCID)

# Directories of Research Data Repositories

- http://www.re3data.org
- http://databib.org/index_subjects.php
- http://oad.simmons.edu/oadwiki/Data_repositories
- http://www.force11.org
- http://www.nature.com/sdata/data-policies/repositories

# Choose a License

Recommended: CC-BY and CC0

**Attribution**
**CC BY**

This license lets others distribute, remix, tweak, and build upon your work, even commercially, as long as they credit you for the original creation. This is the most accommodating of licenses offered. Recommended for maximum dissemination and use of licensed materials.

CC0 (datasets) "No rights reserved"

CC0 can be particularly important for the sharing of data and databases, since it otherwise may be unclear whether highly factual data and databases are restricted by copyright or other rights. Databases may contain facts that, in and of themselves, are not protected by copyright law.

CC0 is recommended for data and databases and is used by hundreds of organizations. It is especially recommended for scientific data. Although CC0 doesn't legally require users of the data to cite the source, it does not take away the moral responsibility to give attribution, as is common in scientific research.

http://creativecommons.org/licenses/

# Simplest Approach

1. Create a public entry for your dataset with a persistent unique identifier
   - Go to zenodo.org, create an account
   - Create an entry for your dataset
2. Specify the metadata
   - Including license -- choose from http://www.creativecommons.org/licenses
3. Upload/point to the data

Voilà!  Figshare will give you a data citation

figshare

creative common

Rv1155,aroG,
icl,Rv1264,thy
Rv0223c,lipJ,Rv1
115    25      cyp130,Rv
20    TB31.7,Rv1264,mscL
1     fabG1,
13    mmaA4,bphD,Rv1264,m
18    TB31.7,cyp130,aroG,
5     pth,ethR,clpP,glbN,
14    pknD,lipJ,fabH,Rv1
10    mmaA4,Rv1264,groE
2     mmaA4,Rv1264,thy
      pepD,Rv1264,thy
      pknD,pepD,fab

# Ideal Approach

1. **Find a repository that your community uses, if there is not one then organize one!**
2. Create a public entry for your dataset with a persistent unique identifier
   - Create an entry for your dataset
3. Specify the metadata required by that repository using metadata standards for that community
   - Including license -- choose from http://www.creativecommons.org/licenses
4. Upload/point to the data
5. Get a data citation from the repository

# What to Show in a GPF

★ Cite each of your datasets like you would cite another paper

★ Citation includes publication date, date of retrieval, repository, and persistent identifier

★ If there is a data paper, cite it

## Data Citation Format

Cite this: Garijo, Daniel; Xie, Lei; Zhang, Yinliang; Gil, Yolanda; Xie, Li; Kinnings, Sarah; Bourne, Phil (2013) Highly connected drug file figshare.
http://dx.doi.org/10.6084/m9.figshare.776887
Retrieved 11:05, Feb 20, 2015 (GMT)

Authors

Date of publication

Time of retrieval

Permanent unique identifier

Name

Repository

# How to Sketch a Workflow

1. Compile the command line invocation to all your codes
   - ★ Input data, parameters, configuration files
   - ★ Include data preparation codes

2. Consider how the data flows from code to code

3. Starting with the input data, work your way to the results

4. If any steps were done with manual intervention, indicate that

5. Create subworkflows if it gets large

Marine Metadata Interoperability

CF MetaData

ISO 19115

WaterML2.0

# Simplest Approach

★ Datasets should have general-purpose metadata specified (creator, date, name, etc.)

# Ideal Approach

★ Dataset characteristics should be explained in detail

★ Domain-specific metadata should be documented

★ Availability of related datasets should be documented

# What to Show in the Paper



The US Long Term Ecological Research Network

★ Mention that the persistent identifier for your data has pointers to its metadata and includes a detailed description of the data

★ Optionally, include the metadata also as supplemental material

★ If there is a data paper, cite it

```
function enEdition(){
    /* Ne rien faire mode edit
    if( encodeURIComponent(documen
turn;
    // /&preload=/

    if ( !wgPageName.match(/Discussion
    var diff = new Array();
    var status; var pecTraduction; var
    var avancementTraduction; var avance

/* *********** Parser *********** */
    var params = document.location.searc
gth).split('&');
    var i = 0;
    var tmp; var name;
    while ( i < params.length )
    {
        tmp = params[i].split('=');
        name = tmp[0];
        switch( name ) {
            case 'status':
                status = tmp[
```

# Simplest Approach

1. Create a public entry for your software with a persistent unique identifier

   - Post on your web site and use a PURL, upload to figshare as you would data and get a DOI

2. Specify the metadata

   - Including license -- choose from http://opensource.org/licenses, preferably Apache v2.0

3. Specify desired citation

# Ideal Approach

1. Learn to use a code repository that allows version tracking and collaborative software development

   - GitHub, BitBucket, etc.

2. Create a public entry for your software with a persistent unique identifier

3. Specify the metadata

   - Including license -- choose from http://opensource.org/licenses, preferably Apache v2.0

4. Specify desired citation

```javascript
function enEdition(){
    /* Ne rien faire mode edit
    if( encodeURIComponent(documen
turn;
    // /&preload=/

    if ( !wgPageName.match(/Discussion
    var diff = new Array();
    var status; var pecTraduction; var
    var avancementTraduction; var avance

/* *********** Parser *********** */
    var params = document.location.searc
gth).split('&');
    var i = 0;
    var tmp; var name;
    while ( i < params.length )
    {
        tmp = params[i].split('=');
        name = tmp[0];
        switch( name ) {
            case 'status':
                status = tmp
```

# Choosing an Open Source License

★ Copyright: automatically applied to software when it is created to grant *the creator* exclusive rights as an intellectual property

★ **Open source license**: reduce constraints and enable software developers to make their source code available to public

1. "Copyleft" license (ex: GNU General Public License (GPL))
2. "Permissive" license (ex: Apache 2 or MIT licenses)

★ **Open Source Initiative**

  ★ Choose a license from: http://opensource.org/licenses
  ★ Recommend that you choose a permissive license
    ★ Apache v2

# What to Show in a GPF

★ Cite each piece of software that you use (preparation, analysis, visualization) like you would cite another paper

  ★ Citation similar to data but includes software version

★ If there is a software paper, cite it

## Software Citation Format

Garijo, Daniel;Xie, Lei; Zhang, Yinliang; Gil, Yolanda; Xie, Li (2013) Tool for computing anomalies, GitHub. V.1 http://dx.doi.org/10.5281/zenodo.18765 Retrieved 11:05, Feb, 15, 2015 (GMT)

**Version**

Authors

Date of publication

Time of retrieval

Permanent unique identifier

Name

Repository

# Simplest Approach

1. Describe as much metadata as you can in your software site

   1. Document the basic metadata

   2. If you use a code repository, there is some basic structure you can follow

# Ideal Approach

1. Use software registry
   - http://www.ontosoft.org/portal, csdms.colorado.edu, etc.
   - Guides through questions to provide metadata
2. Save the metadata as HTML, XML,…
3. Post the metadata on your code site

# What to Show in the Paper

★ Mention that the persistent identifier location for your software points to its metadata

★ Optionally, include the software metadata as supplemental material

★ If there is a software paper, cite it

## PIHM [Christopher Duffy]

### Identify

#### Locate - Unique description

**What is the software called ?**

○ PIHM

**What is a short description for this software ?**

○ PIHM is a multiprocess, multi-scale hydrologic model where the major hydrological processes are fully coupled using the semi-discrete finite volume method. PIHM is a physical model for surface and groundwater, "tightly-coupled" to a GIS interface. PIHMgis which is open source, platform independent and extensible. The tight coupling between GIS and the model is achieved by developing a shared data-model and hydrologic-model data structure.

Initial metadata was retrieved from http://csdms.colorado.edu/wiki/Model:PIHM

**What are general categories (keywords, labels) for this software ?**

○ Hydrology
○ Basins
○ Continental

**Is there a project website for the software ?**

○ http://www.pihm.psu.edu/pihm_home.html

### Understand

#### Trust - Quality and ratings

**Who created this software? (Project, Organization, Person, Initiative, etc.)**

○ Christopher Duffy

**Are there any additional contributors of note for this software ?**

○ Mukesh Kumar
○ Gopal Bhatt

# Simplest Approach

1. Describe the workflow in text
   - Data + software + workflow
   - Specify unique identifiers for data and software, versions, credit all sources
2. Develop a workflow sketch
   - Capture high-level dataflow across components
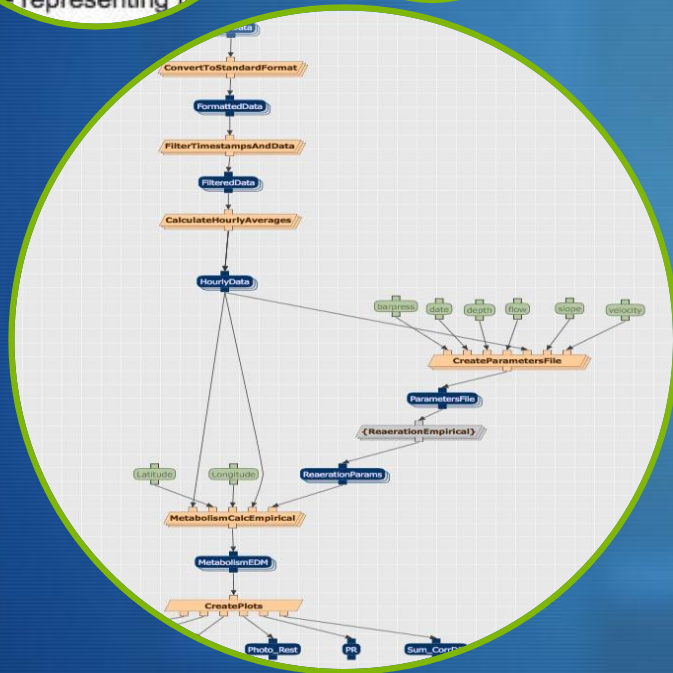3. For provenance, include a summary or an execution trace

**5** Provenance documentation

**6** Methods documentation

# Ideal Approach

1. Describe the workflow in text
   - Data + software + workflow
   - Specify unique identifiers for data and software, versions, credit all sources
2. Develop a workflow sketch
   - Capture high-level dataflow across components
3. Specify the formal workflow using a workflow system, electronic notebook, etc.
   - Command lines + parameter values
   - Dataflow across components
4. Include the provenance record
   - If generating it automatically, preferably using a standard (e.g., PROV)
5. Publish the workflow and provenance record in a publicly accessible repository (eg figshare, myExperiment, etc)
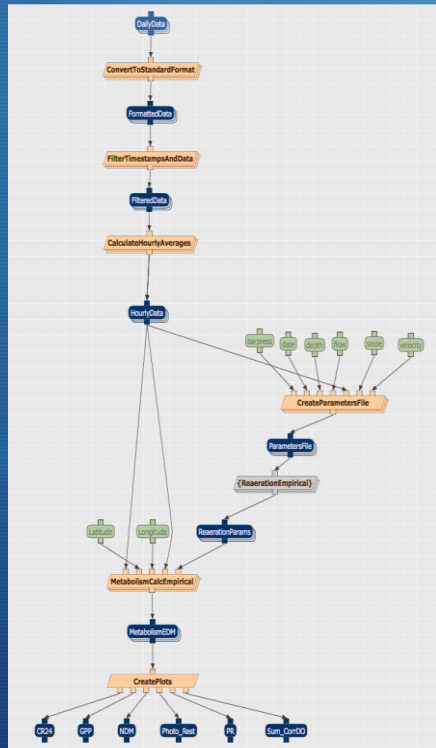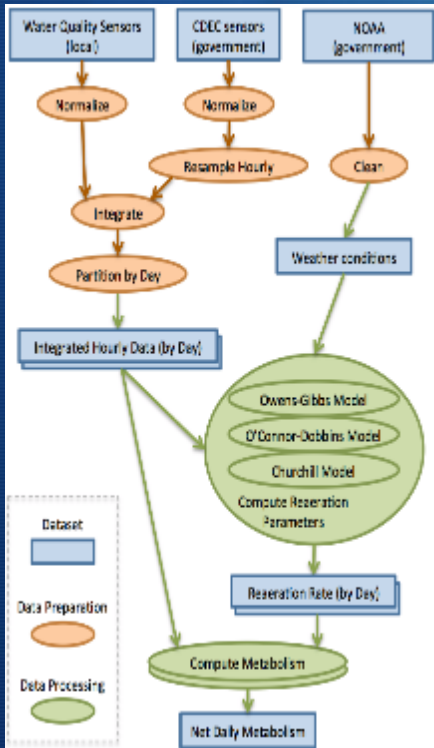6. Get a unique persistent identifier for the workflow, the provenance, or both

**5** Provenance documentation

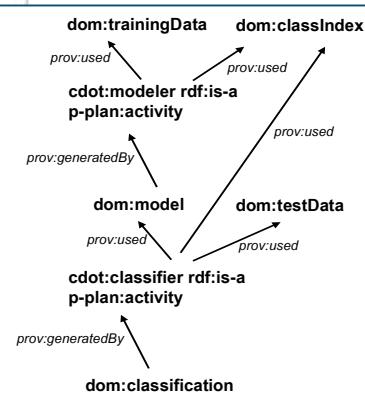**6** Methods documentation

# What to Show in the Paper

★ Describe workflow in text and provide a workflow sketch

  ★ Optionally, provide the formal workflow or lab notebook, use a persistent identifier, and cite it

★ Include a summary of the execution traces as supplementary material, or use a persistent identifier and cite it

  ★ Optionally, include instead the provenance records using a standard like W3C PROV

dom:trainingData    dom:classIndex

*prov:used*                    *prov:used*

cdot:modeler rdf:is-a
p-plan:activity

*prov:generatedBy*                        *prov:used*

dom:model    dom:testData

*prov:used*         *prov:used*

cdot:classifier rdf:is-a
p-plan:activity

*prov:generatedBy*

dom:classification

**# Entities**
ex:testData1 a prov:Entity .
ex:model1 a prov:Entity .
ex:classification1 a prov:Entity .

**# Activities**
ex:Classifier1 a prov:Activity .

**# Usage and Generation relations between entities and activities**
ex:Classifier1
    prov:used ex:testData1 ;
    prov:used ex:model1 .

ex:classification1
    prov:wasGeneratedBy
        ex:Classifier1 .
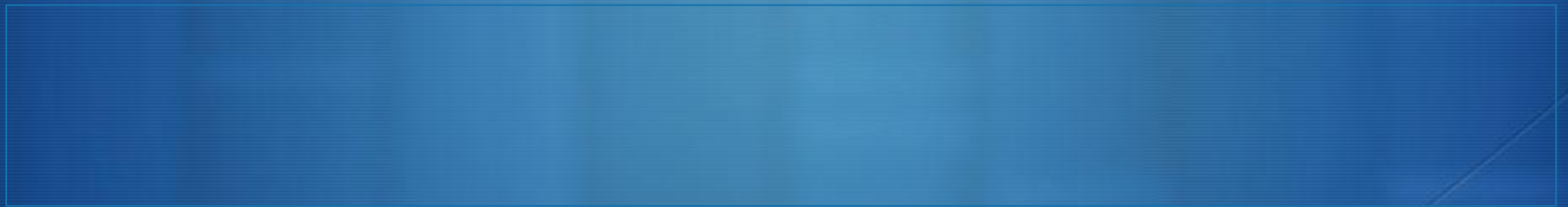
# What to Show in the Paper

★ Authors have a persistent unique identifier

  ★ Use www.orcid.org

ORCiD

# EXTRA SLIDES: Workflows from [Zhu, Ridley, and Deng 2016]

Research paper

# Simulating electron and ion temperature in a global ionosphere thermosphere model: Validation and modeling an idealized substorm

Jie Zhu [a,*], Aaron J. Ridley [a], Yue Deng [b]

[a] Department of Atmospheric, Oceanic and Space Sciences, University of Michigan, Ann Arbor, MI, USA
[b] Department of Physics, University of Texas, Arlington, USA
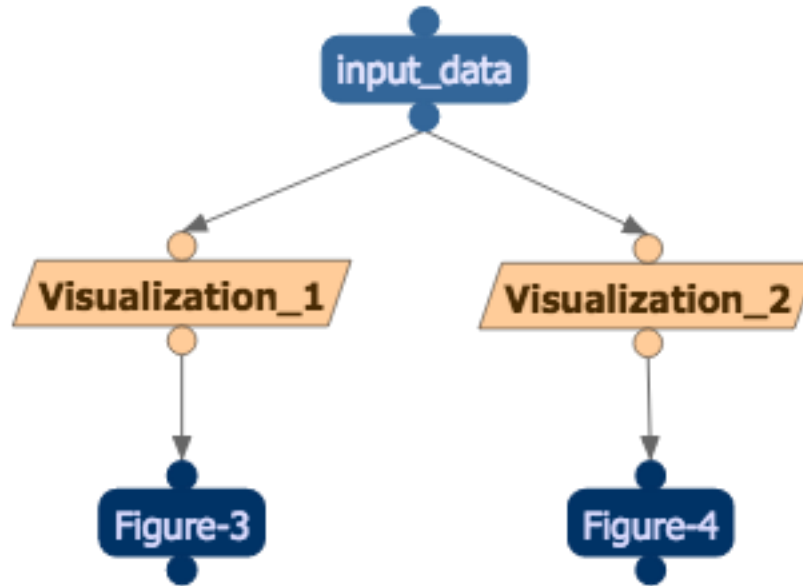
## A R T I C L E   I N F O

## A B S T R A C T

Electron and ion temperatures control many chemical and physical processes in the ... mosphere system. Recently, improved electron and ion energy equations were i... Global Ionosphere Thermosphere Model (GITM). The source energy of the electro... includes thermal conduction, heating due to photoionization, elastic collisions wi... inelastic collisions with neutrals, auroral precipitation, and heat flux from inner ... source terms in the ion temperature ($T_i$) equation include thermal conduction, and e... electrons and neutrals. The new implementation of $T_e$ improved the ionospheric de... high latitudes with respect to IRI. The improved GITM also reproduced the diurnal v... observed by incoherent scatter radars at low and middle latitudes. The model was us... idealized substorm statistically described by Clausen et al. (2014). It was found that ...

# Figures 5 & 6 & 9

# Figure 7
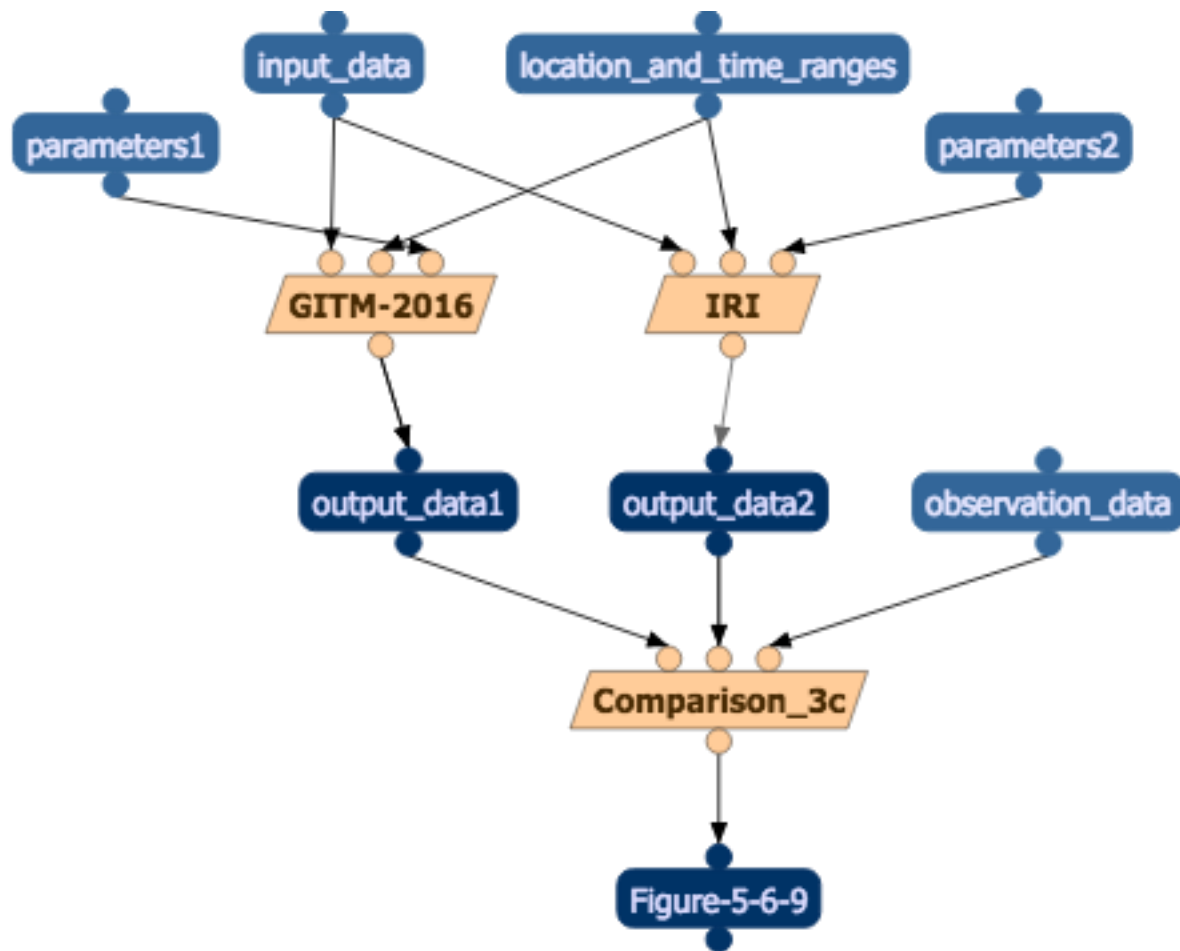
# Figure 8