

The *human* factor in big-data and eco-informatics



Ankur R Desai
University of Wisconsin-Madison
ESA 2018, SYMP 9
7 Aug 2018
New Orleans, LA

Acknowledgments

- Jack Williams, UW-Madison
- Kim Novick, University of Indiana
- Michael Dietze, Boston University
- Stefan Metzger and Wendy Gram, NEON/Battelle
- Kathleen Weathers, Cary Institute
- Ben Bond Lamberty, DOE PNNL
- Enablsh FAIR Data Project
- And many contributors...
- + Support from NSF BIO (DEB, EF/MSB, ABI, RCN) and AGS, DOE TES Ameriflux, Battelle/NEON, NASA Carbon Cycle, NOAA, USGCRP

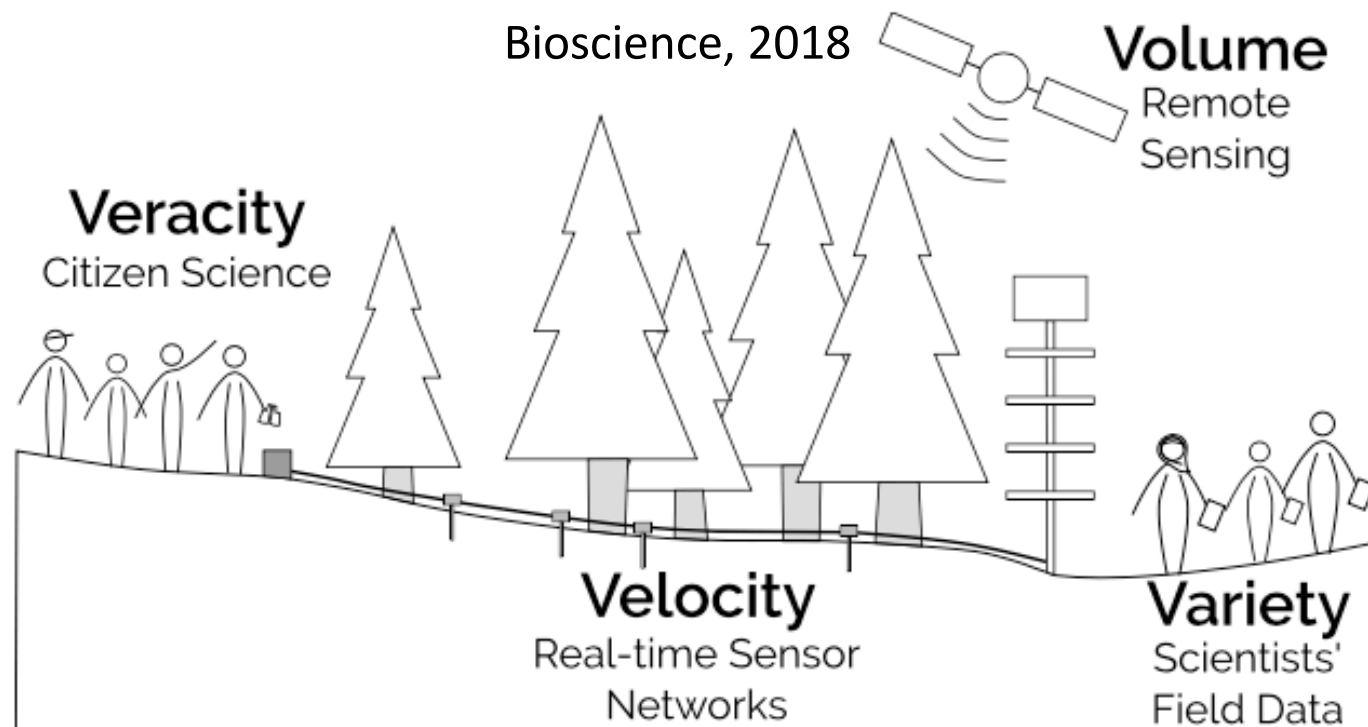
Take Homes

- Big data is not just about data volume
 - Data/code diversity, accessibility, and metadata matter
- Tackling challenges in informatics is a key to solving the scientific reproducibility crisis
 - Big data is really about the people, ethics, networks
- Ecologists are well-positioned to be a leader here
 - If we invest in **people** and **infrastructure**

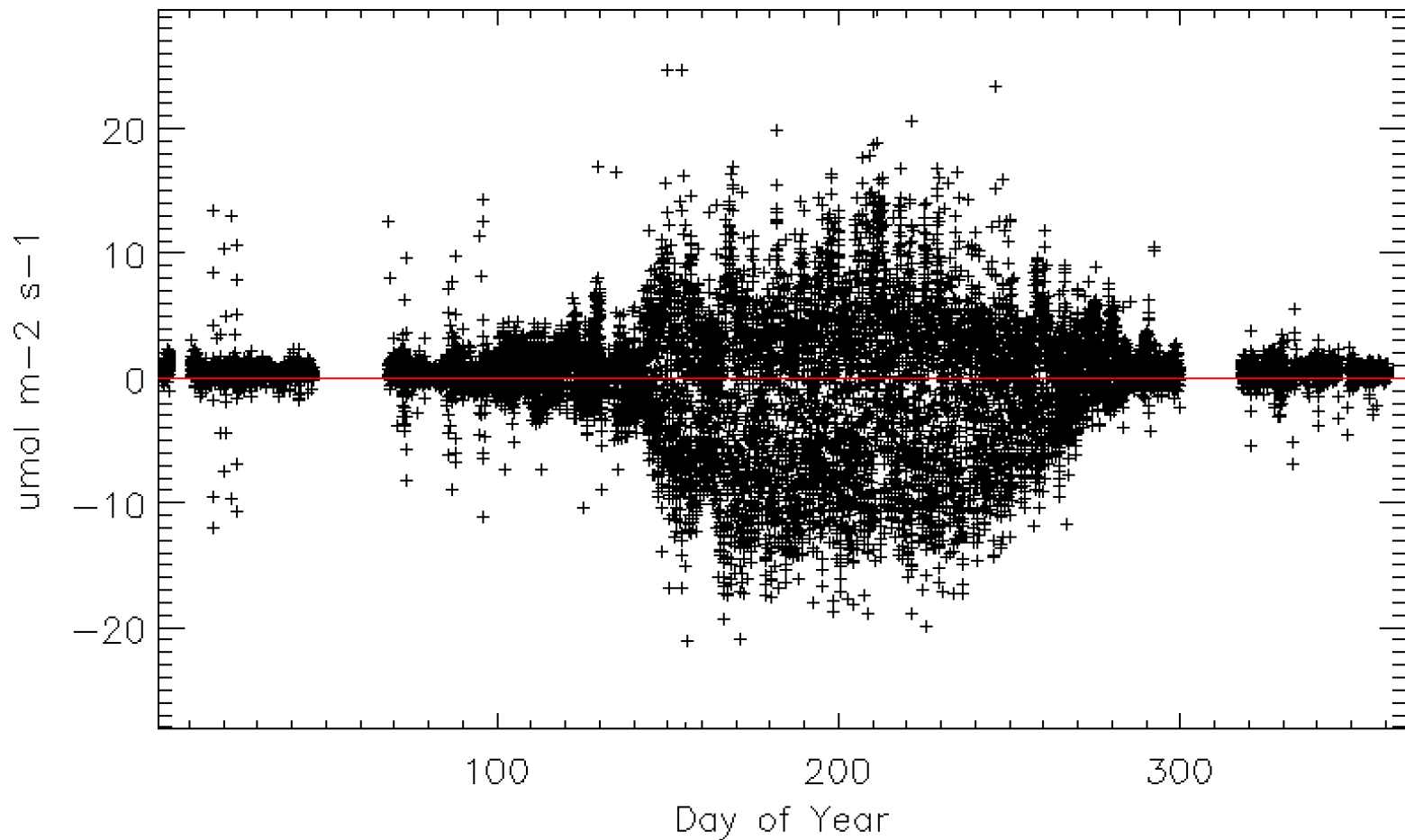
You *all* already do
“big-data”
“fusion”
“forecasting”
and “informatics”

Situating Ecology as a Big-Data Science: Current Advances, Challenges, and Solutions

SCOTT S. FARLEY, ANDRIA DAWSON, SIMON J. GORING AND JOHN W. WILLIAMS



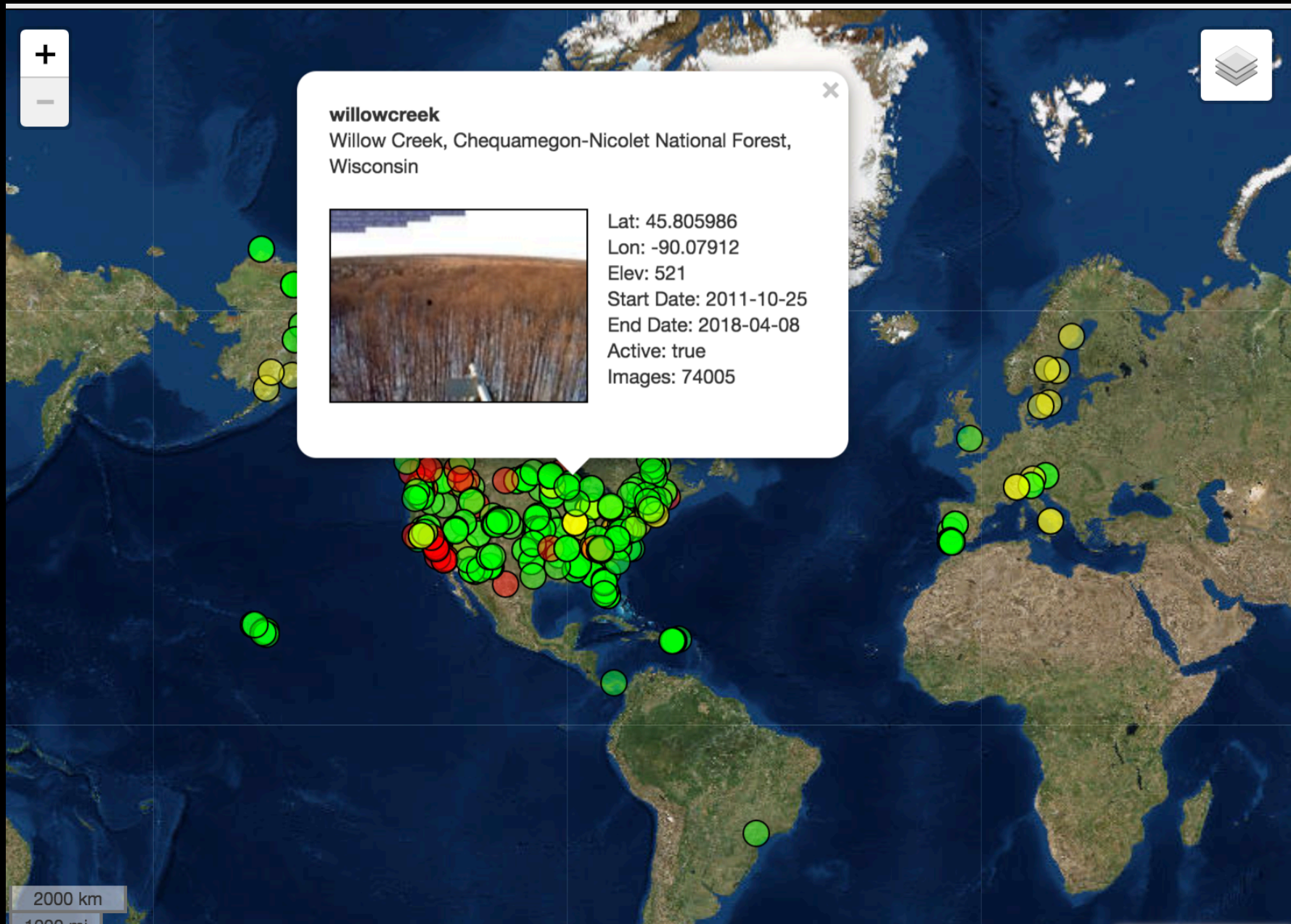
DATA!!! Om nom nom...







<https://phenocam.sr.unh.edu/webcam/gallery/>

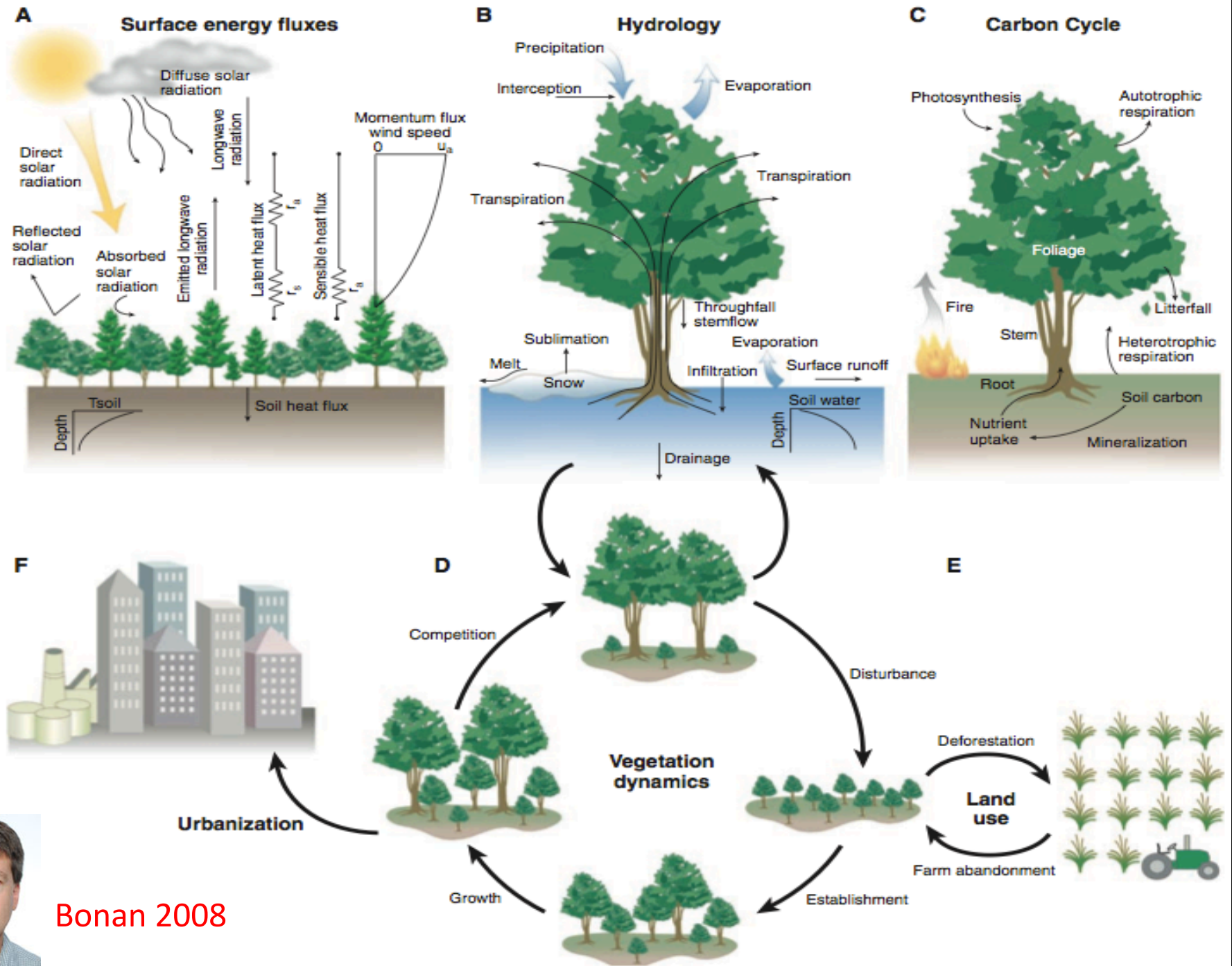


Andrew Richardson, Northern Arizona U

What are these data good for?

- Understand, measure, and predict the fate of global-warming greenhouse gases and how that influences ongoing and future climate change
 - Atmospheric and ecological theories of vegetation-climate **feedbacks**
 - Long-term, **multi-scale** observations of soil and vegetation carbon and water use
 - **Fusing** these to confront numerical models of land surface biophysics, ecosystem dynamics, and atmospheric forcing/feedbacks

Forests in Flux

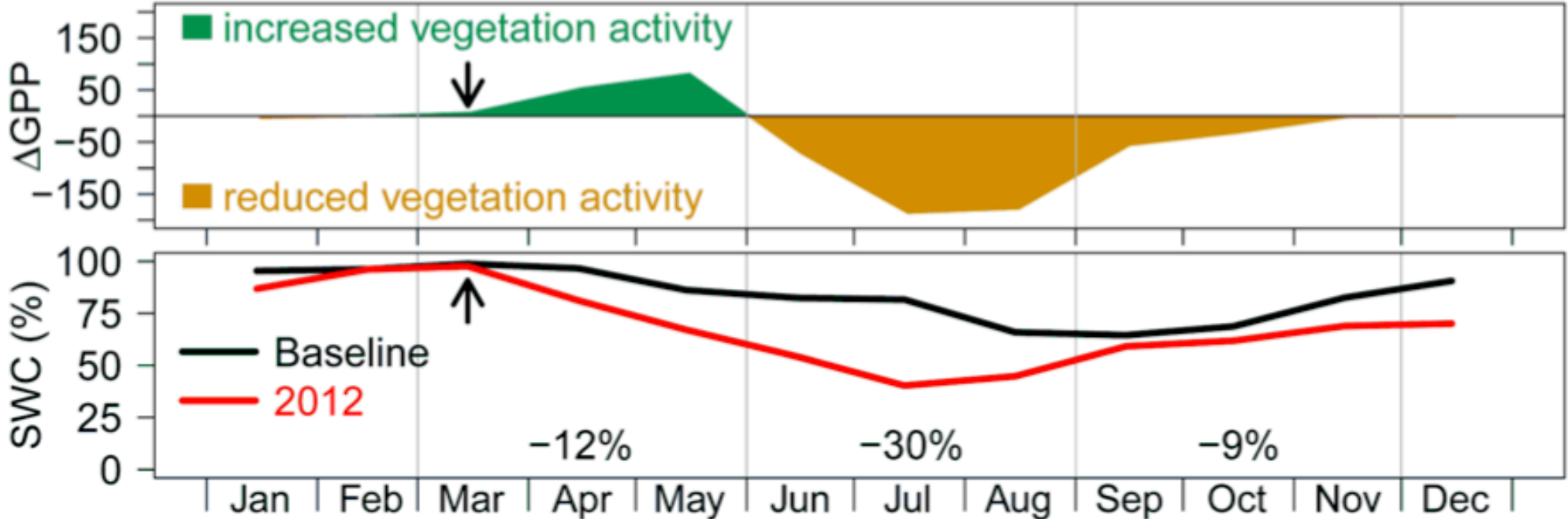
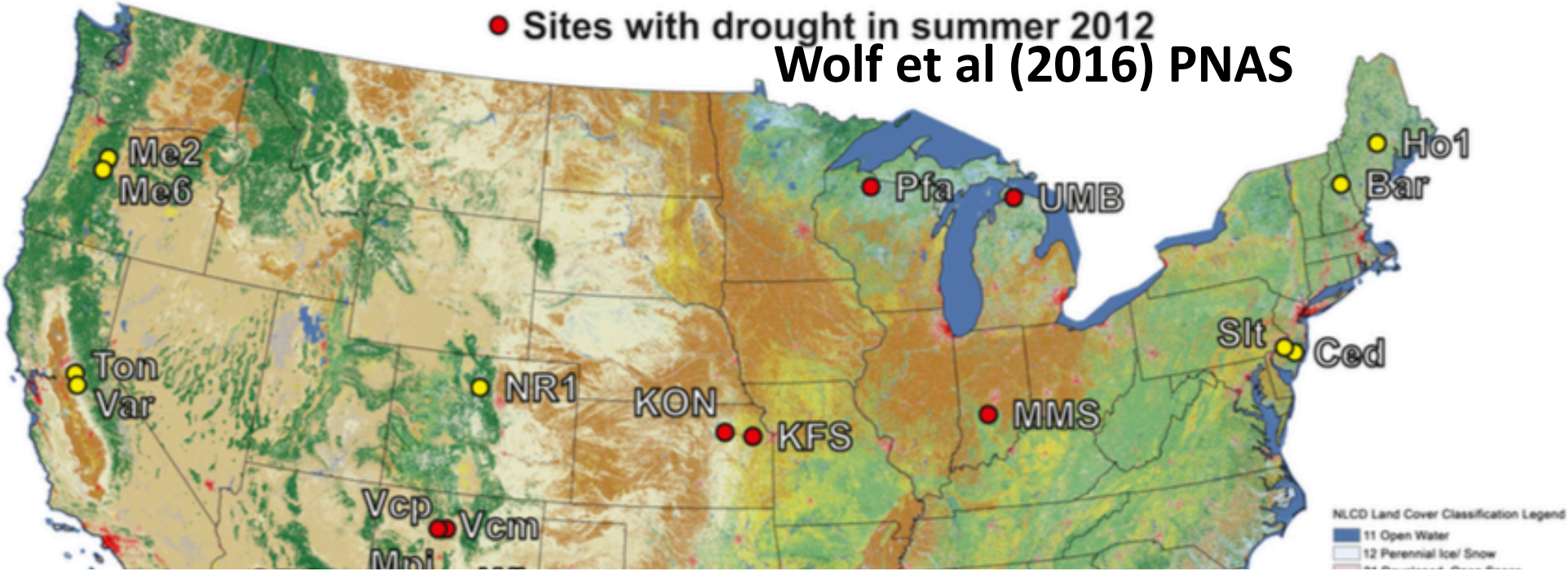


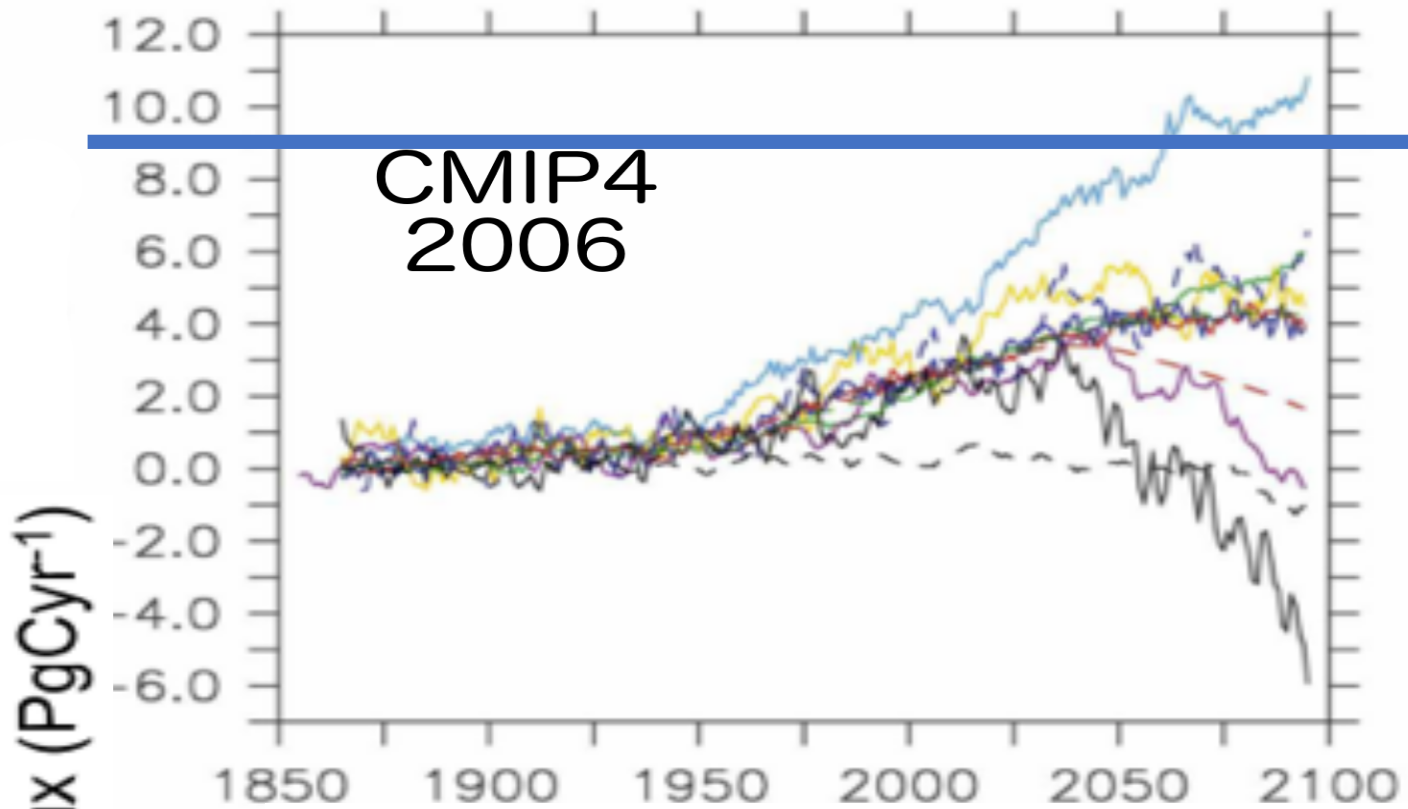
Bonan 2008



Wolf et al (2016) PNAS

● Sites with drought in summer 2012



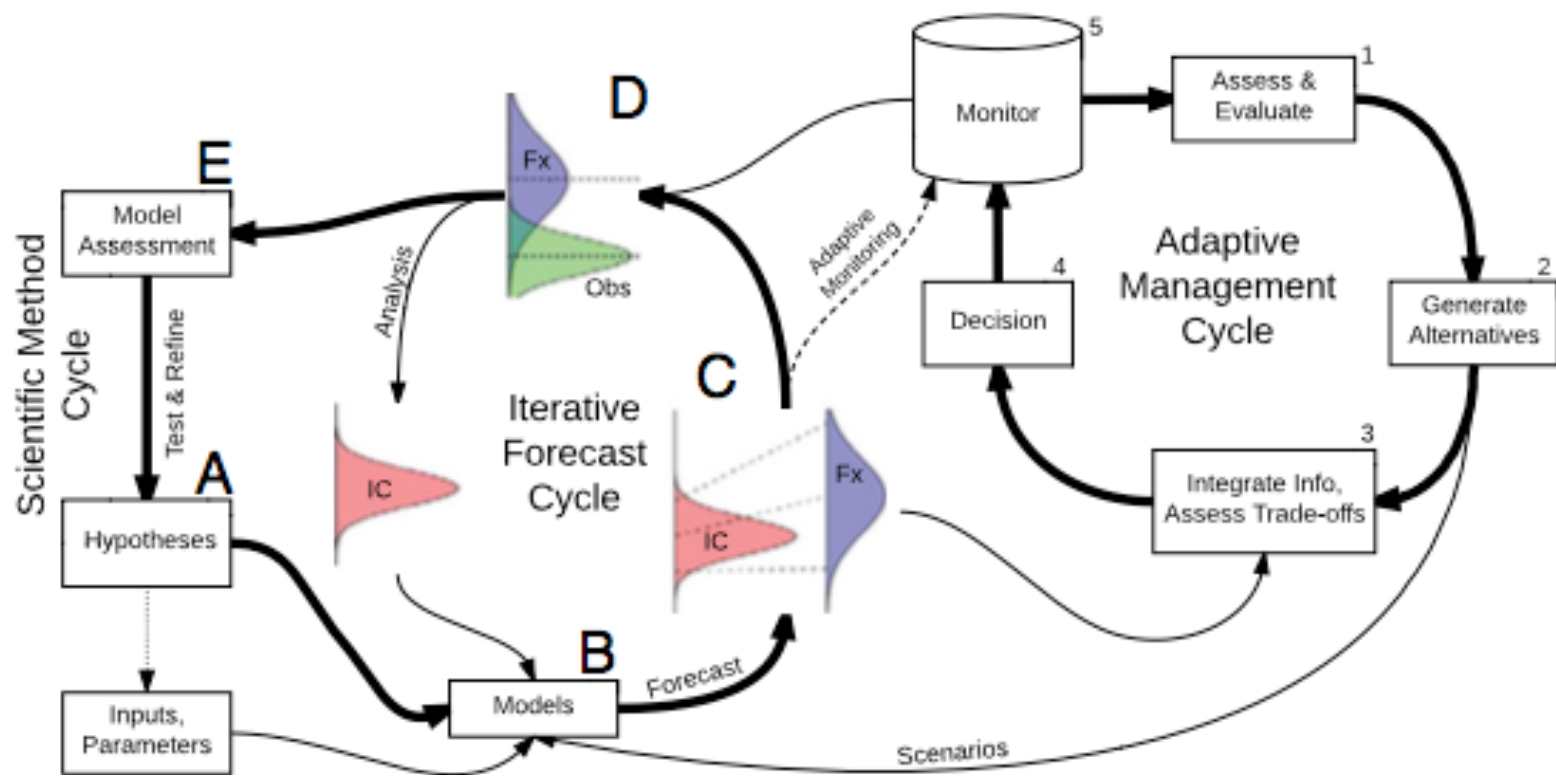


Annual land flux (PgCyr⁻¹)

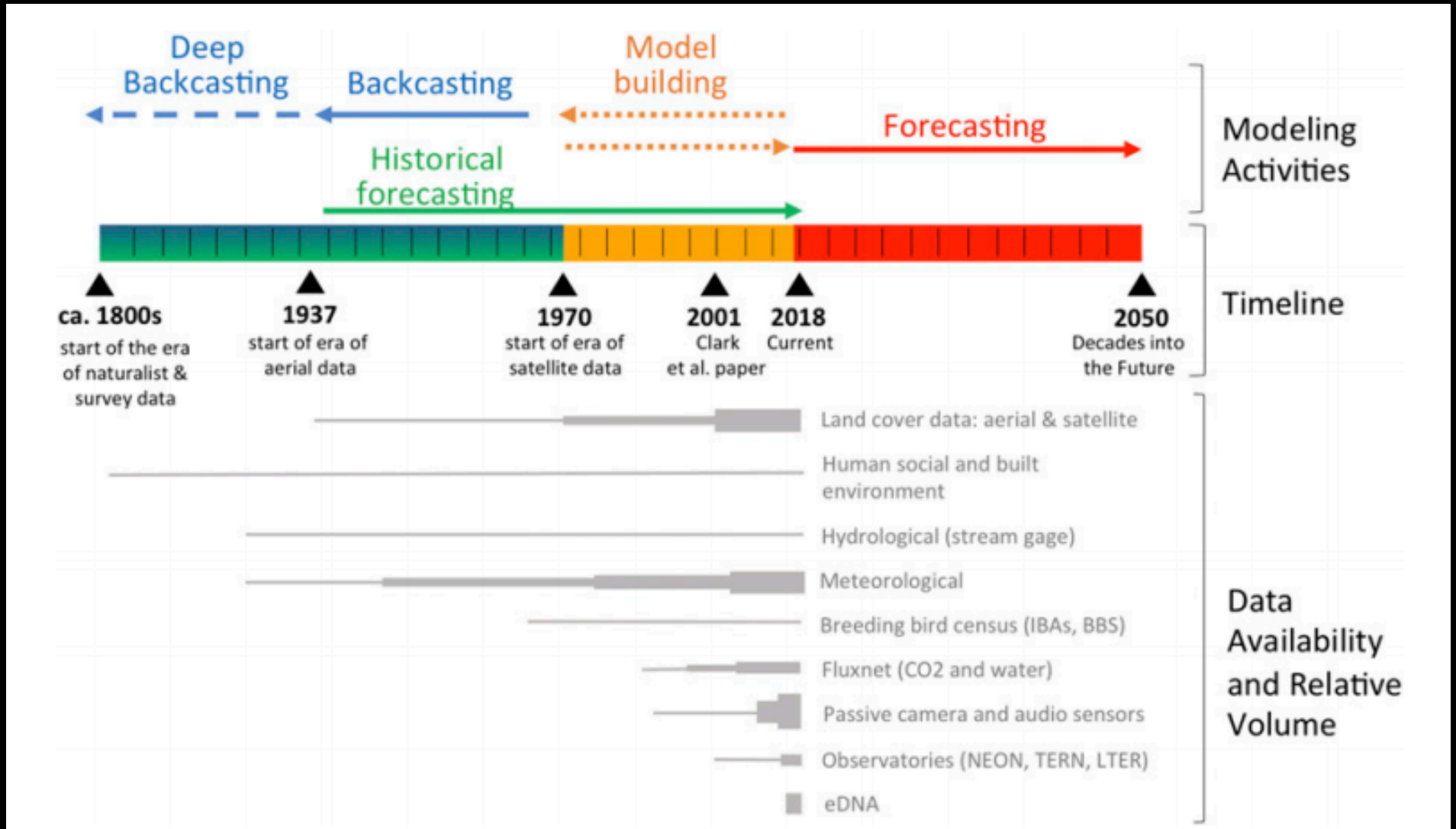
10
8
6
4
2
0
-2
-4
-6
-8

Iterative near-term ecological forecasting: Needs, opportunities, and challenges

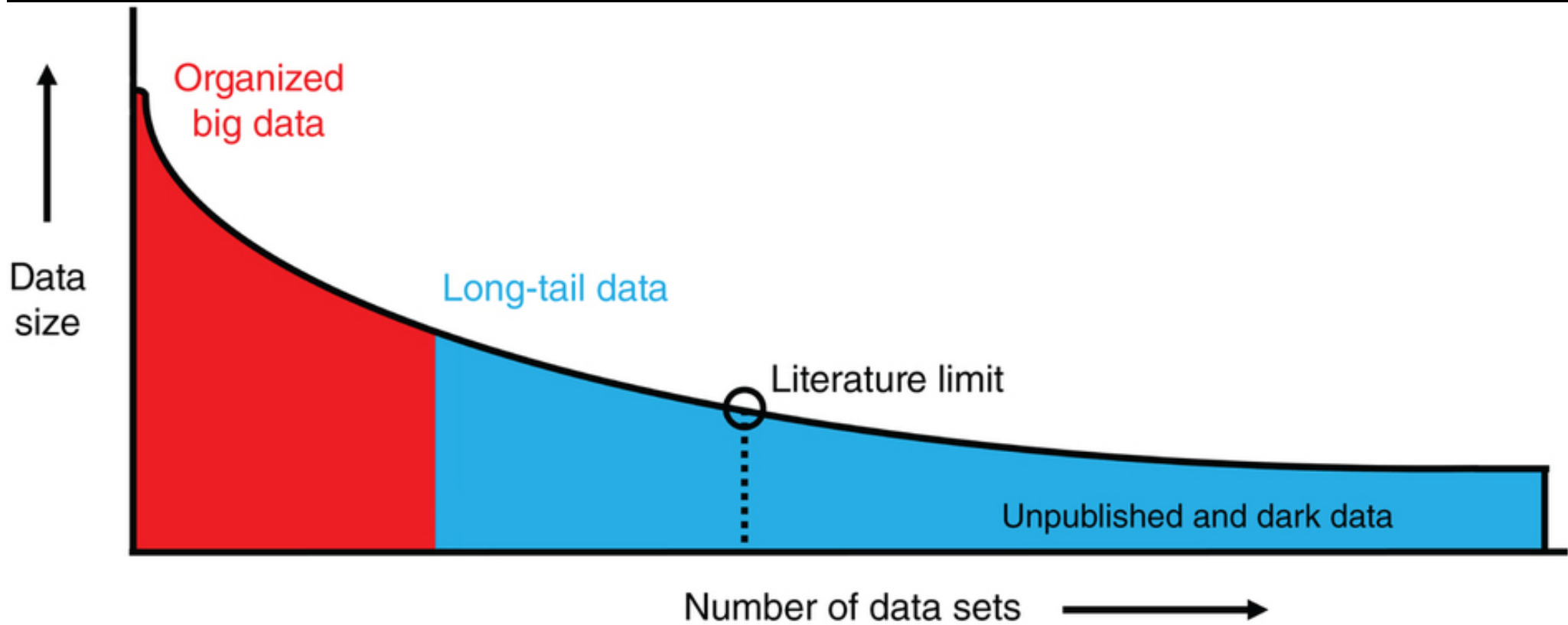
Michael C. Dietze^{a,1}, Andrew Fox^b, Lindsay M. Beck-Johnson^c, Julio L. Betancourt^d, Mevin B. Hooten^{e,f,g}, Catherine S. Jarnevich^h, Timothy H. Keittⁱ, Melissa A. Kenney^j, Christine M. Laney^k, Laurel G. Larsen^l, Henry W. Loescher^{k,m}, Claire K. Lunch^k, Bryan C. Pijanowskiⁿ, James T. Randerson^o, Emily K. Read^p, Andrew T. Tredennick^{q,r}, Rodrigo Vargas^s, Kathleen C. Weathers^t, and Ethan P. White^{u,v,w}



Observations are big and long!



And hard to extract from literature!



Ferguson et al., 2014
Nature Neuroscience



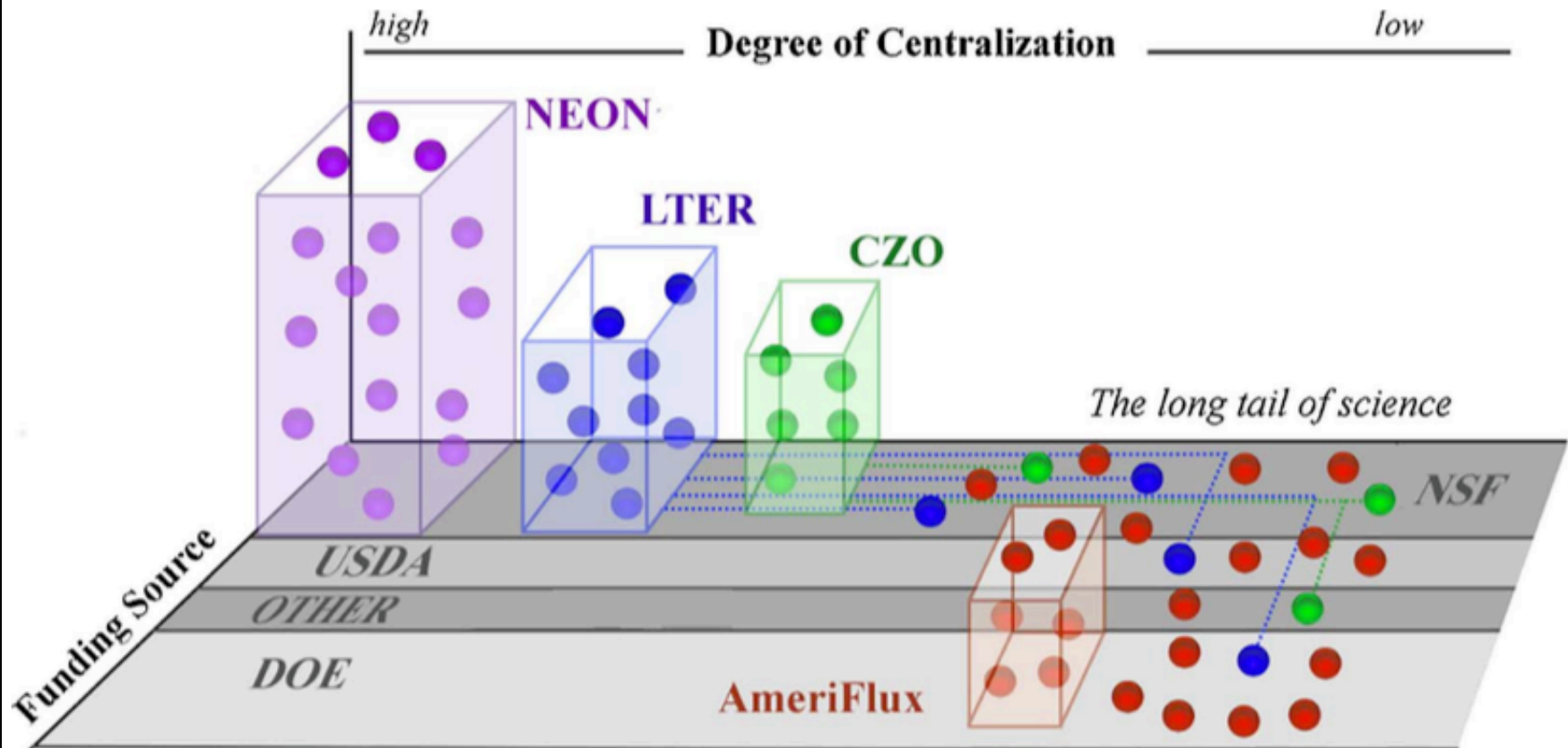
Contents lists available at ScienceDirect

Agricultural and Forest Meteorology

journal homepage: www.elsevier.com/locate/agrformet

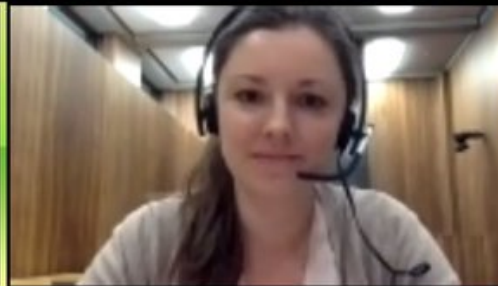
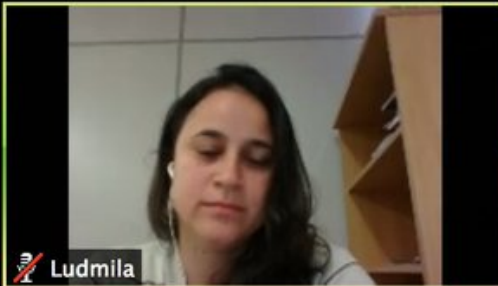
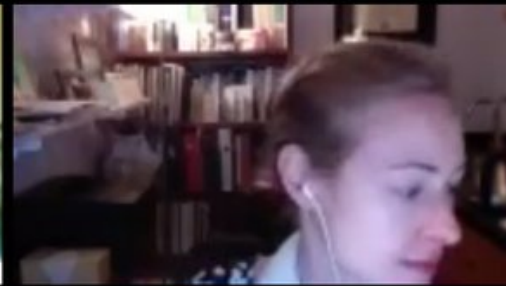
The AmeriFlux network: A coalition of the willing

K.A. Novick^{a,*}, J.A. Biederman^b, A.R. Desai^c, M.E. Litvak^d, D.J.P. Moore^e, R.L. Scott^b, M.S. Torn^f







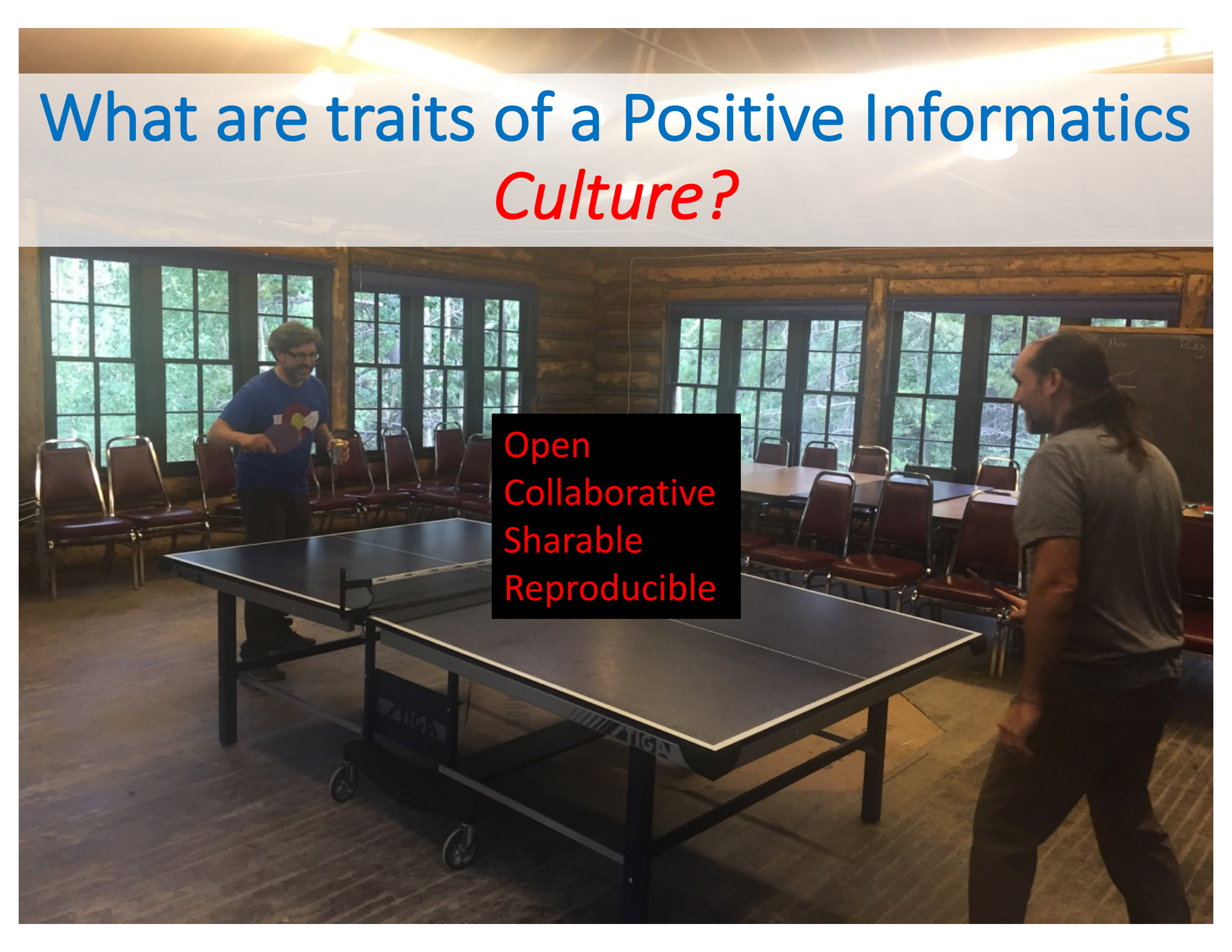


robynsmyth

uli

Giovanna Flaim

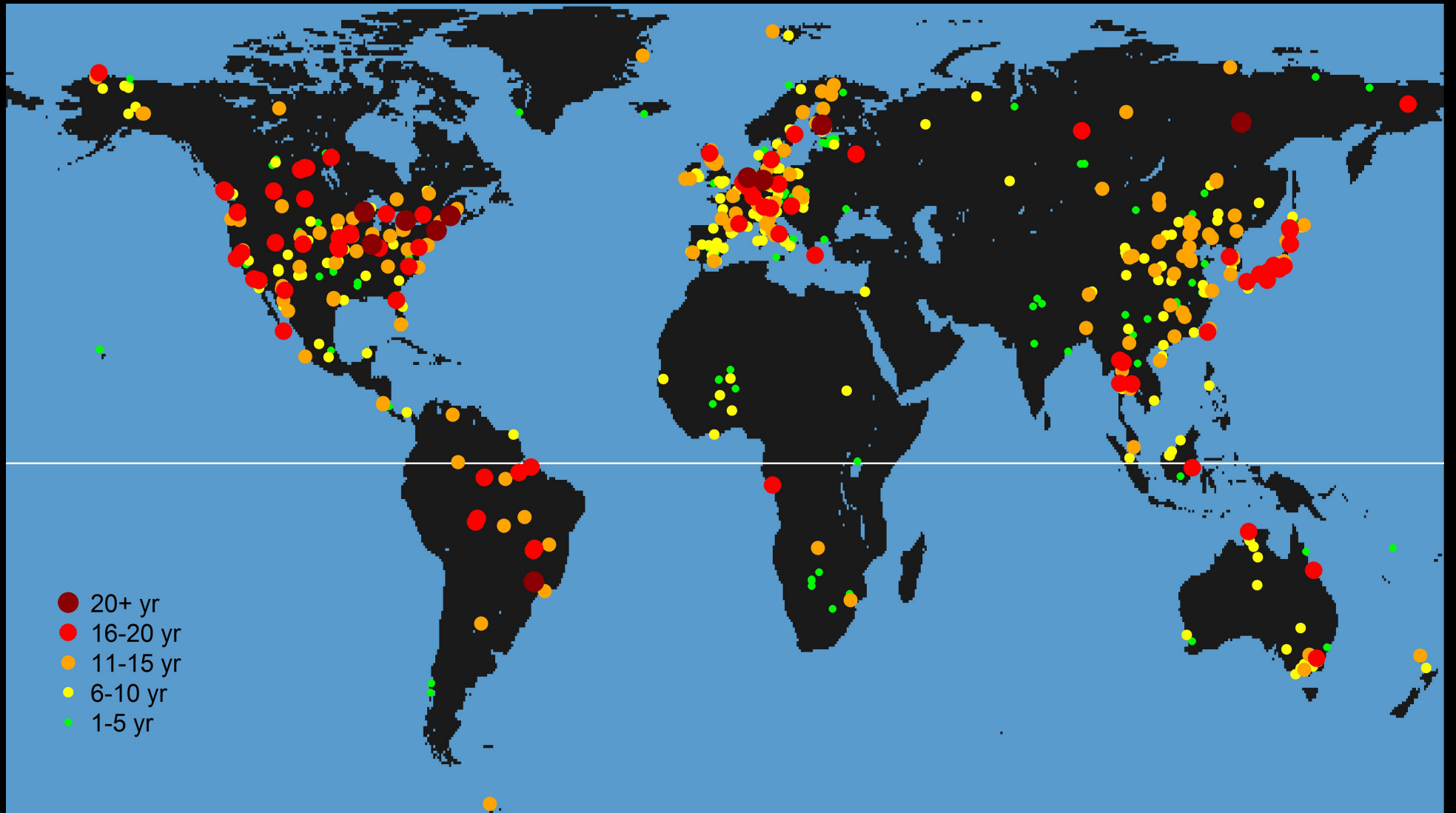
What are traits of a Positive Informatics *Culture?*

A photograph of two men playing ping pong in a log cabin. The man on the left is wearing a blue t-shirt and glasses, and is in the middle of a stroke. The man on the right is wearing a grey t-shirt and has a ponytail. The room has large windows, wooden walls, and several tables and chairs. A blackboard is visible in the background.

Open
Collaborative
Sharable
Reproducible

OPEN







The AmeriFlux network: A coalition of the willing

K.A. Novick^{a,*}, J.A. Biederman^b, A.R. Desai^c, M.E. Litvak^d, D.J.P. Moore^e, R.L. Scott^b, M.S. Torn^f

Table 3

Strengths and weaknesses of AmeriFlux's bottom-up approach.

| Feature of Approach | Associated Strengths | Associated Weaknesses |
|---|---|---|
| Voluntary, PI-driven research; inclusive approach to network participation | Diverse research questions; interdisciplinarity; strong sense of community Good spatial and temporal representativeness of many biome types. | Lack of incentives for data sharing. Insecurity of funding for many sites. Underrepresentation of some biomes. |
| Lack of standardization of instrumentation and processing “collaborative” data policy | Flexibility in methodological approach can advance observation theory. Promotes cross-disciplinary perspectives; strengthens interpersonal connections within the network; promotes incentive for PIs to submit data | Biases related to instrument design and processing can challenge cross-site syntheses. Large, multi-author papers are sometimes challenging to write, presenting a disincentive for network end-users. |
| Network oriented around a relatively few core observations (i.e. fluxes and meteorological drivers) | Few required variables makes it easier for sites to join the network | Inconsistent submission of non-biometeorological data across sites, which when present provides important ecological context for the fluxes, and guides model development. |

FLUX COURSE

#fluxcourse

Annual Summer Course in Flux Measurements and Advanced Modeling:
Training scientists in observations and models to advance carbon cycle science

David JP Moore¹, Kim Novick²

¹University of Arizona, School of Natural Resources and the Environment, ²School of Public and Environmental Affairs, Indiana University - Bloomington



WED COS 82

Dave Moore

Kim Novick

The all singing
all dancing
ecologist?

Forming
communities
of practice by
cross training
graduate
students in
empirical and
modelling
approaches

3:20 pm

Rm 245

The fluxcourse has run for 10 years. You can find instructor videos and lessons from past courses at our Gitbook (QR Code) <https://tinyurl.com/FluxcourseGitbook>



www.fluxcourse.org

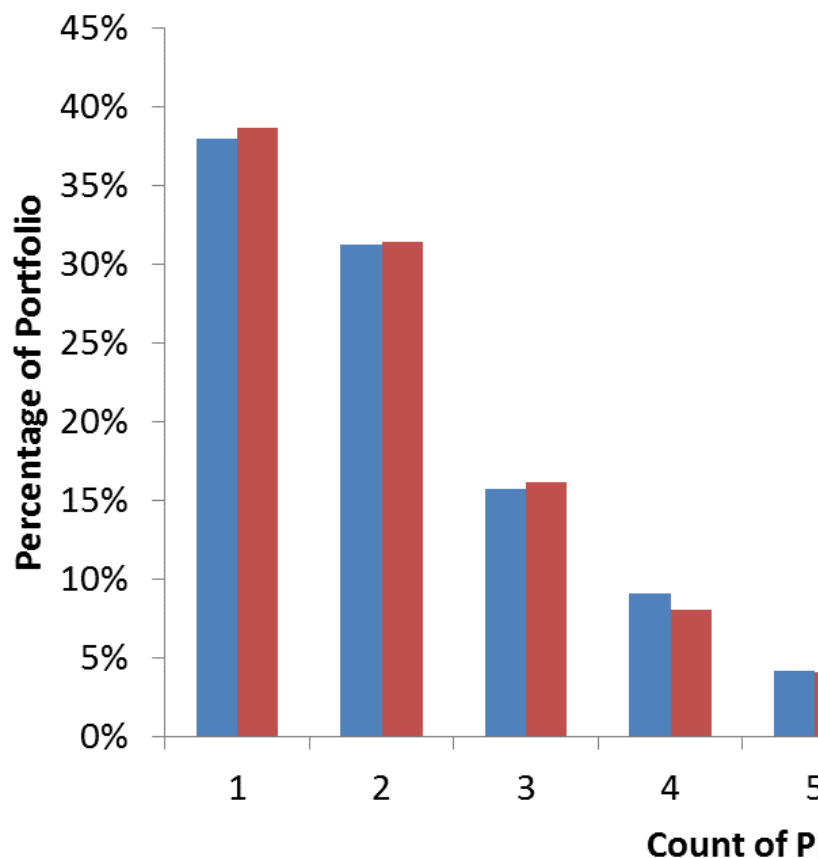


COLLABORATIVE



Most scientific pro multi-PI, multi-ins

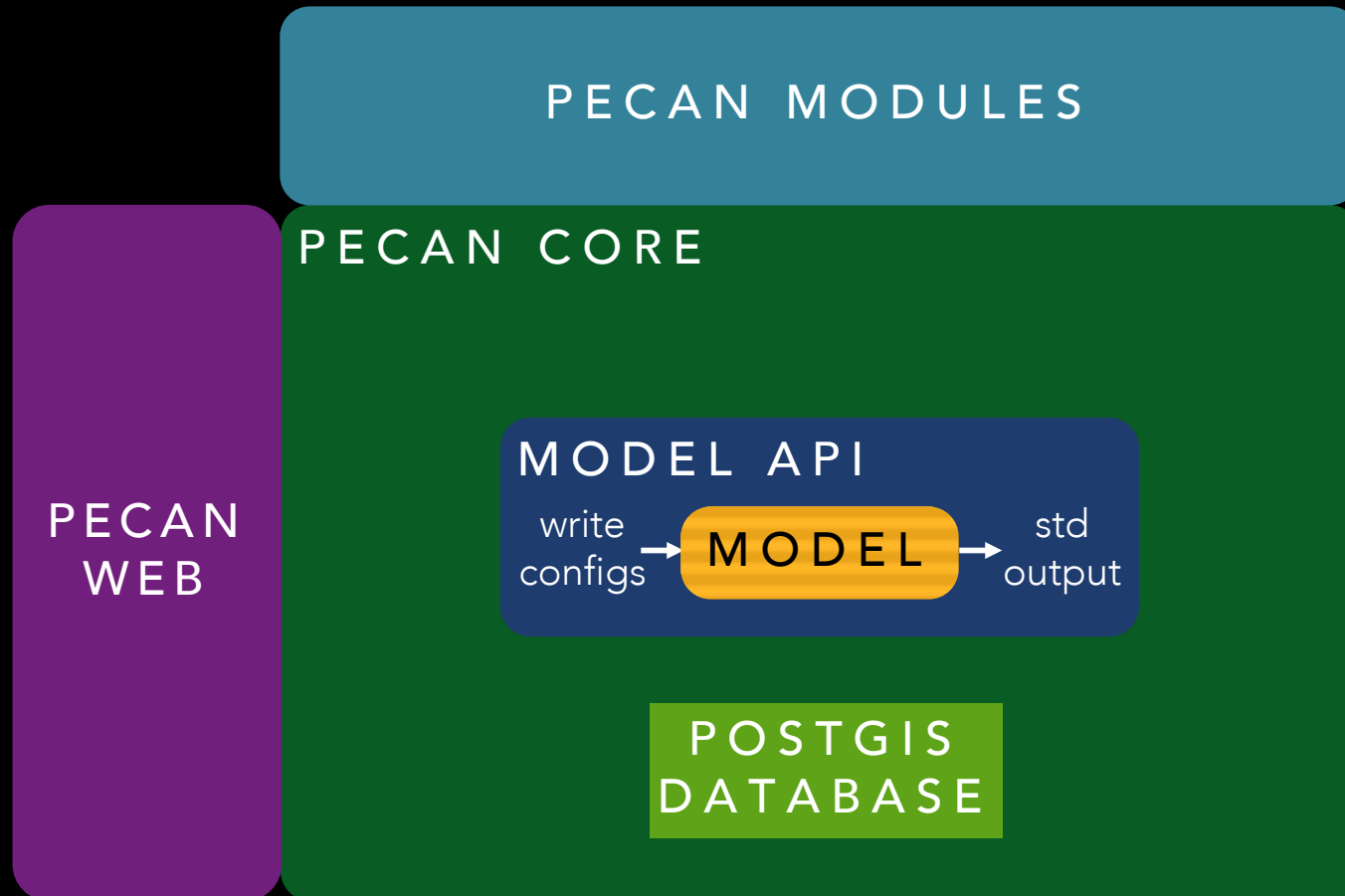
2009-2012 DEB Core Pro



<https://deblog.nsfbio.com/2013/07/23/deb-numbers-award-size-and-duration/>

Climate control of terrestrial carbon exchange across biomes and continents

Chuixiang Yi¹, Daniel Ricciuto², Runze Li³, John Wolbeck¹, Xiyan Xu¹, Mats Nilsson⁴, Luis Aires^{5,117}, John D Albertson^{6,117}, Christof Ammann^{7,117}, M Altaf Arain^{8,117}, Alessandro C de Araujo^{9,117}, Marc Aubinet^{10,117}, Mika Aurela^{11,117}, Zoltán Barcza^{12,117}, Alan Barr^{13,117}, Paul Berbigier^{14,117}, Jason Beringer^{15,117}, Christian Bernhofer^{16,117}, Andrew T Black^{17,117}, Paul V Bolstad^{18,117}, Fred C Bosveld^{19,117}, Mark S J Broadmeadow^{20,117}, Nina Buchmann^{21,117}, Sean P Burns^{22,117}, Pierre Cellier^{23,117}, Jingming Chen^{24,117}, Jiquan Chen^{25,117}, Philippe Ciais^{26,117}, Robert Clement^{27,117}, Bruce D Cook^{28,117}, Peter S Curtis^{29,117}, D Bryan Dalrymple^{30,117}, Ebba Dellwik^{31,117}, Nicolas Delpechier^{32,117}, Ankur R Desai^{33,117}, Sabina Dore^{34,117}, Danilo Dragoni^{35,117}, Bert G Drake^{36,117}, Eric Dufrêne^{32,117}, Allison Dunn^{37,117}, Jan Elbers^{38,117}, Werner Eugster^{21,117}, Matthias Falk^{39,117}, Christian Feigenwinter^{40,117}, Lawrence B Flanagan^{41,117}, Thomas Foken^{42,117}, John Frank^{43,117}, Juerg Fuhrer^{7,117}, Damiano Gianelle^{44,117}, Allen Goldstein^{45,117}, Mike Goulden^{46,117}, Andre Granier^{47,117}, Thomas Grünwald^{48,117}, Lianhong Gu^{2,117}, Haiqiang Guo^{49,117}, Albin Hammerle^{50,117}, Shijie Han^{51,117}, Niall P Hanan^{52,117}, László Haszpra^{53,117}, Bernard Heinesch^{10,117}, Carole Helfter^{54,117}, Dimmie Hendriks^{55,117}, Lindsay B Hutley^{56,117}, Andreas Ibrom^{57,117}, Cor Jacobs^{38,117}, Torbjörn Johansson^{58,117}, Marjan Jongen^{59,117}, Gabriel Katul^{60,117}, Gerard Kiely^{61,117}, Katja Klumpp^{62,117}, Alexander Knohl^{21,117}, Thomas Kolb^{34,117}, Werner L Kutsch^{63,117}, Peter Lafleur^{64,117}, Tuomas Laurila^{11,117}, Ray Leuning^{65,117}, Anders Lindroth^{58,117}, Heping Liu^{66,117}, Benjamin Loubet^{23,117}, Giovanni Manca^{67,117}, Michal Marek^{68,117}, Hank A Margolis^{69,117}, Timothy A Martin^{70,117}, William J Massman^{43,117}, Roser Matamala^{71,117}, Giorgio Matteucci^{72,117}, Harry McCaughey^{73,117}, Lutz Merbold^{74,117}, Tilden Meyers^{75,117}, Mirco Migliavacca^{76,117}, Franco Miglietta^{77,117}, Laurent Misson^{78,117,118}, Meelis Mölder^{58,117}, John Moncrieff^{79,117}, Russell K Monson^{79,117}, Leonardo Montagnani^{80,81,117}, Mario Montes-Helu^{84,117}, Eddy Moors^{82,117}, Christine Moureaux^{10,83,117}, Mukufute M Mukelabai^{84,117}, J William Munger^{85,117}, May Myklebust^{65,117}, Zoltán Nagy^{86,117}, Asko Noormets^{87,117}, Walter Oechel^{88,117}, Ram Oren^{89,117}, Stephen G Pallardy^{90,117}, Kyaw Tha Paw U^{39,117}, João S Pereira^{59,117}, Kim Pilegaard^{57,117}, Krisztina Pintér^{86,117}, Casimiro Pio^{91,117}, Gabriel Pita^{92,117}, Thomas L Powell^{93,117}, Serge Rambal^{94,117}, James T Randerson^{46,117}, Celso von Randow^{95,117}, Corinna Rebmann^{64,117}, Janne Rinne^{96,117}, Federica Rossi^{77,117}, Nigel Roulet^{97,117}, Ronald J Ryel^{98,117}, Jorgen Sagerfors^{4,117}, Nobuko Saigusa^{99,117}, María José Sanz^{100,117}, Giuseppe-Scarascia Mugnozza^{101,117}, Hans Peter Schmid^{102,117}, Guenther Seufert^{103,117}, Mario Siqueira^{89,117}, Jean-François Soussana^{62,117}, Gregory Starr^{104,117}, Mark A Sutton^{105,117}, John Tenhunen^{106,117}, Zoltán Tuba^{86,117,118}, Juha-Pekka Tuovinen^{11,117}, Riccardo Valentini^{107,117}, Christoph S Vogel^{108,117}, Jingxin Wang^{109,117}, Shaoqiang Wang^{110,117}, Weiguo Wang^{111,117}, Lisa R Welp^{112,117}, Xuefa Wen^{110,117}, Sonia Wharton^{113,117}, Matthew Wilkinson^{20,117}, Christopher A Williams^{114,117},



Standardized inputs and outputs

Provenance: Transparent & Repeatable

Accessible interface

Reusable tools for execution, analysis, visualization

No central repository!



For code or data!

SHARABLE



Sharing is caring...

- The National Ecological Observatory Network is a \$450 million NSF set of coordinated U.S. ecological observing sites to address grand challenges in global change
 - The “space telescope” of ecology
- Community resource – consistent instruments on all sites, open data, documentation for every variable REST/JSON API for access
- But can this infrastructure support ecology?

eddy-covariance usability tools: eddy4R-Docker image

- Docker = shipping container system for code



Geosci. Model Dev. Discuss., doi:10.5194/gmd-2016-318, 2017

Manuscript under review for journal Geosci. Model Dev.

Published: 1 February 2017

© Author(s) 2017. CC-BY 3.0 License.



Geoscientific
Model Development
Discussions



1 **eddy4R: A community-extensible processing, analysis and**
2 **modeling framework for eddy-covariance data based on R,**
3 **Git, Docker and HDF5**

4

5 **Stefan Metzger¹, David Durden¹, Cove Sturtevant¹, Hongyan Luo¹, Natchaya**
6 **Pingintha-Durden¹, Torsten Sachs², Andrei Serafimovich², Jörg Hartmann³,**
7 **Jiahong Li⁴, Ke Xu⁵, Ankur R. Desai⁵**

Enabling FAIR Data

- Findable
 - Accessible
 - Interoperable
 - Re-Usable
-
- <https://www.force11.org/group/fairgroup/fairprinciples>
 - <http://www.copdess.org/enabling-fair-data-project/commitment-to-enabling-fair-data-in-the-earth-space-and-environmental-sciences/>

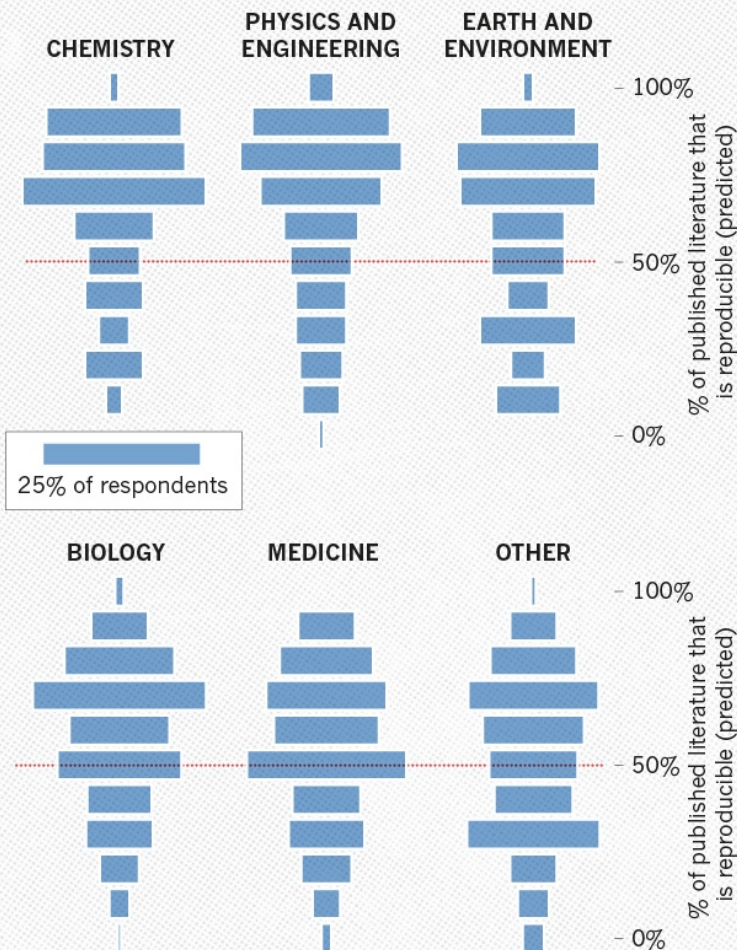
REPRODUCIBLE



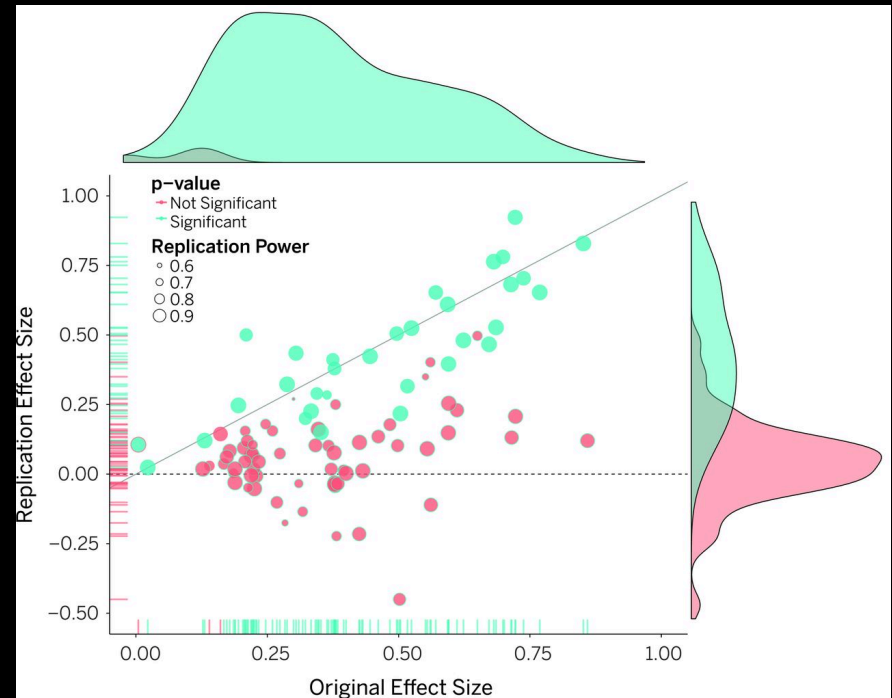
We have a reproducibility crisis...

HOW MUCH PUBLISHED WORK IN YOUR FIELD IS REPRODUCIBLE?

Physicists and chemists were most confident in the literature.



Number of respondents from each discipline:
 Biology 703, Chemistry 106, Earth and environmental 95,
 Medicine 203, Physics and engineering 236, Other 233 ©nature



Estimating the reproducibility of psychological science

Open Science Collaboration^{*†}

* See all authors and affiliations

Science 28 Aug 2015;
 Vol. 349, Issue 6251,
 DOI: 10.1126/science.aac4716

<http://www.nature.com/news/1-500-scientists-lift-the-lid-on-reproducibility-1.19970>

Lack of full metadata is an issue

- Protocol
- Code
- Data
- Filtering and tests
- Experiments and vignettes

- #openexperiment

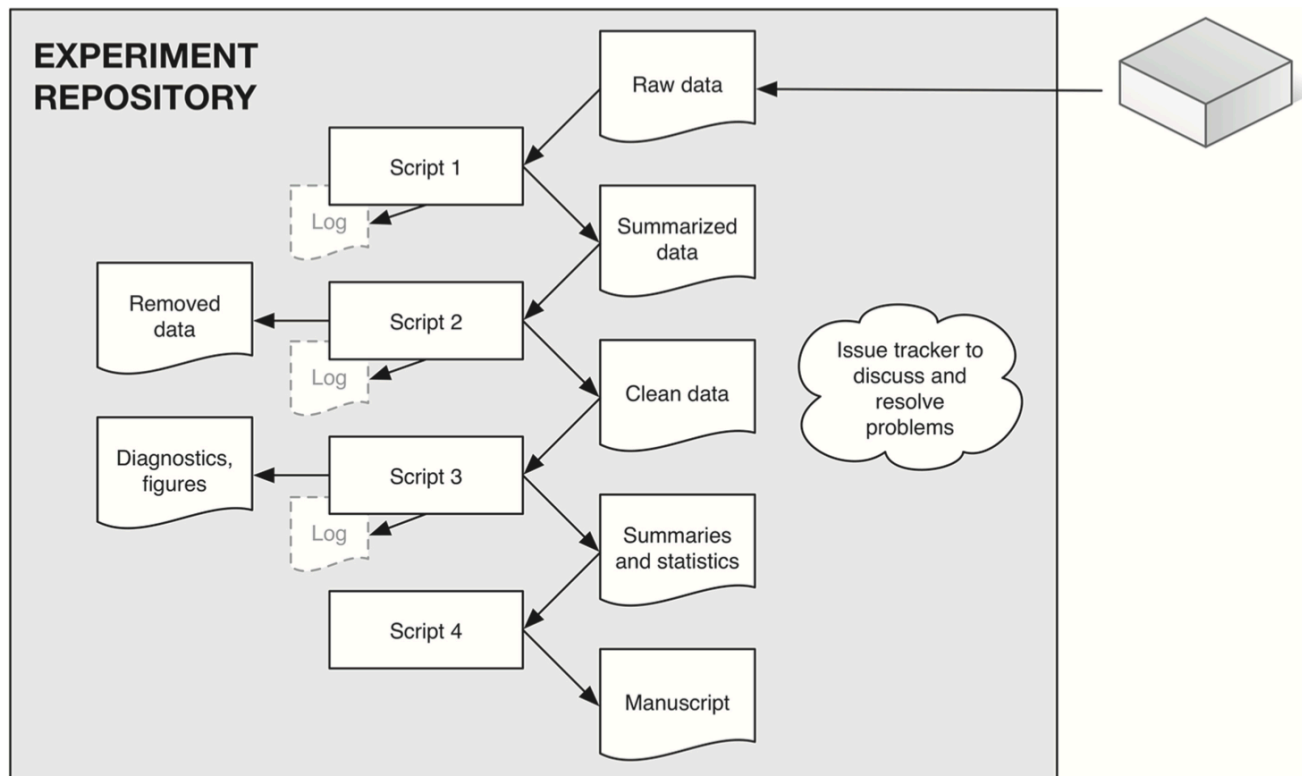
Environmental Research Letters

https://github.com/bpbond/cpcrw_incubation

LETTER

Running an open experiment: transparency and reproducibility in soil and ecosystem science

Ben Bond-Lamberty¹, A Peyton Smith² and Vanessa Bailey²



Missing/problematic data

Number of reps by date and core



↻ Paul Stoy Retweeted



Ben Bond-Lamberty @BenBondLamberty · Aug 1

Terrible internet here in rural Ontario, but: our @Nature paper published today entirely #openscience #opendata. All results & figures reproducible, and you can look through commits to see initial idea, final tweaks, and all my curse words in between.

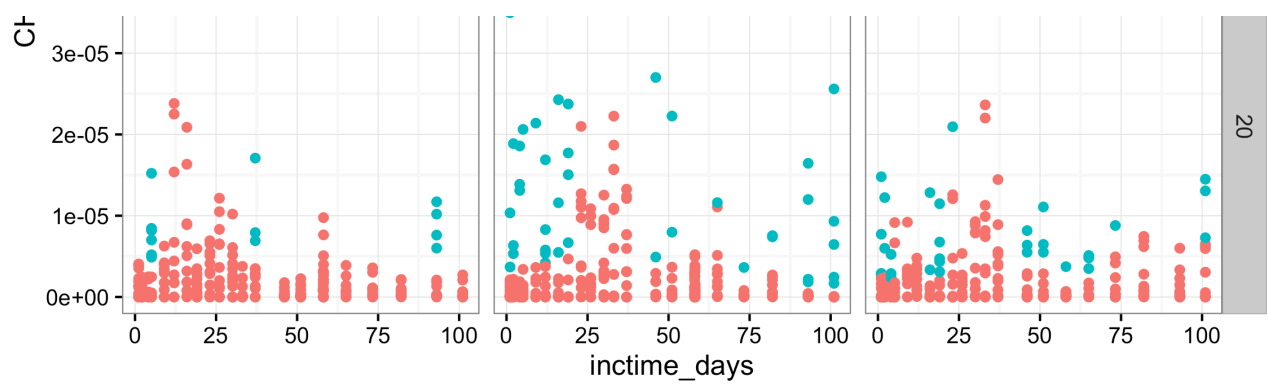


bpbond/rh-changes

Testing for changes in global Rh. Contribute to rh-changes development by creating an account on GitHub.

github.com

CH4_outlier
 ● FALSE
 ● TRUE



Big data is not open, collaborative sharable, nor reproducible if...

- Code to generate/analyze is not reusable by others
 - Github, Docker, DevOps cycle is key to making “big science” happen
- Data lack open, common APIs to access by machines
 - THREDDS, JSON/XML
- Data formats are non-standard, not machine-readable
 - NetCDF, Unidata CF convention as an example in meteorology
 - Ecological Metadata Language (EML)
- Data requires complex authentication methods to access or repositories don't have multiple points of entries, distributed nodes
 - Kill the password!
- Data/code sharing policies limit what you can do
 - Important to set this out by community, be open to ideas beyond intended use
- Data quicklooks, comparisons, documentation on variable names, time steps, units are not easy to find
 - Simple tables, online, vignettes, forums/chat rooms

How do we encourage and support an *open, collaborative, sharable, reproducible* informatics culture?

- Training for students *and* us old farts
- Best practices: GLEON, LTER, Ameriflux
- Science of Team Science
 - NSF RCN MSB Grassroots Global Network Sciences
- Pilot projects for new collaboration methods
 - Will some please invent an actually usable collaborative video-conference platform?
- Funding support for data archival and informatics
 - Digitization/generation of metadata for long-tail data
 - The mantra does NOT have to be centralization
 - CyVerse
- ... what else?



FIELD GUIDE

L. Michelle Bennett
Howard Gadlin
Christophe Marchand



TRUST

It is almost impossible to have trust. Trust provides the foundation for the nearly impossible to succeed.



TEAM EVOLUTION AND DYNAMICS

Research teams form and develop through critical interactions. Their highest potential (Forming, Storming, Norming, Performing) is sustained and further strengthened by positive team dynamics, enabling it to achieve successful outcomes.



VISION

A strong and captivating vision provides a foundation for collaboration. It provides a focal point for team members to coalesce.



COMMUNICATION

Effective communication within and outside a research team contributes to effective group functioning. It depends on a safe environment where team members can openly discuss new scientific ideas and take research in new or unconsidered directions as well as ensure that decisions can take place.



SELF-AWARENESS

Emotional Intelligence is the effective functioning of people. Greater self-awareness improves the quality of team awareness.



RECOGNITION AND SHARING SUCCESS

Individual contributions should be recognized, rewarded in the context of a collaboration. Recognition of all team members should be done thoughtfully in the context of the team and the institution.



LEADERSHIP

Strong collaborative leadership leverages team members' strengths and abilities. Leadership can be derived from the formal leader(s).



CONFLICT AND DISAGREEMENT

Conflict can be both a resource and a challenge—because disagreement can expand thinking, add to a complex scientific problem, and stimulate new ideas for research. A challenge because if it is not handled, conflict impedes effective team functioning and advancement.



MENTORING

Mentoring is an indispensable part of team development. A mentor recognizes the areas in which newer scientists need help and can help coach people. Through mentoring, the development and strengthening of team dynamics.

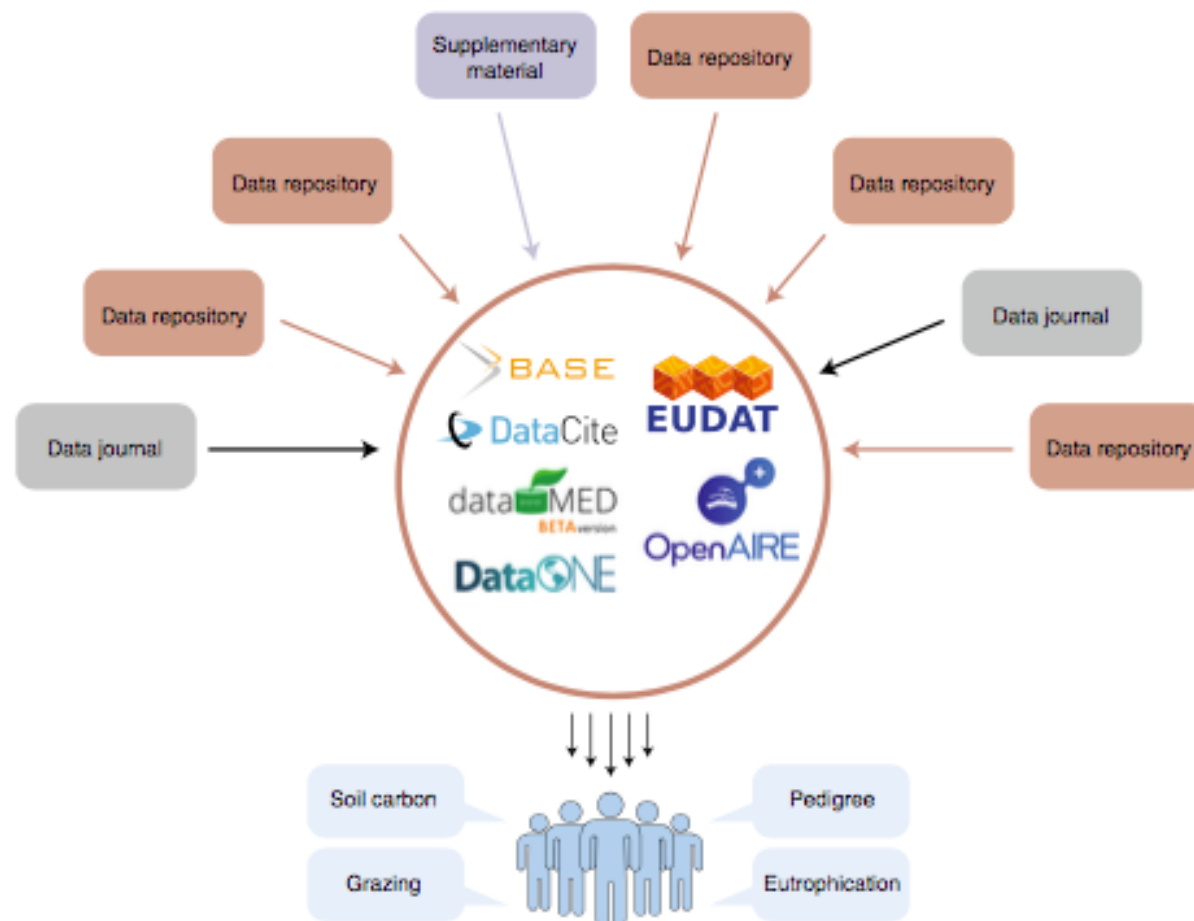


NAVIGATING AND LEVERAGING NETWORKS AND SYSTEMS

Highly collaborative teams can transcend different organizational structures, extending their reach across and beyond their own organization. They often function within the context of a larger network or system.

Navigating the unfolding open data landscape in ecology and evolution

Antica Culina^{1*}, Miriam Baglioni², Tom W. Crowther^{1,3}, Marcel E. Visser¹,
Saskia Woutersen-Windhouver¹ and Paolo Manghi²



THANK YOU!

Ankur Desai, desai@aos.wisc.edu, <http://flux.aos.wisc.edu>

Funding: NSF Advances in Biological Informatics (**ABI-1457897** , **ABI-1062205**)

NSF MSB RCN (**EF-1702991**), NSF NTL LTER (**DEB- 1440297**), DOE Ameriflux Project

