# Cyberinfrastructure: Opportunities and Challenges
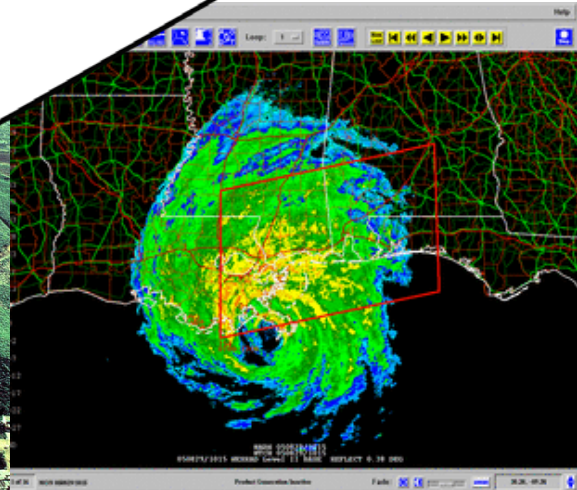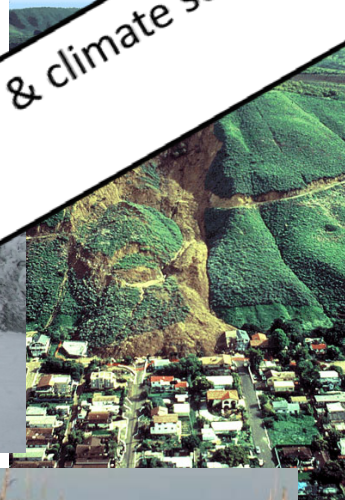
## EarthCube Mesoscale Modeling Workshop

17 December2012

Boulder, CO

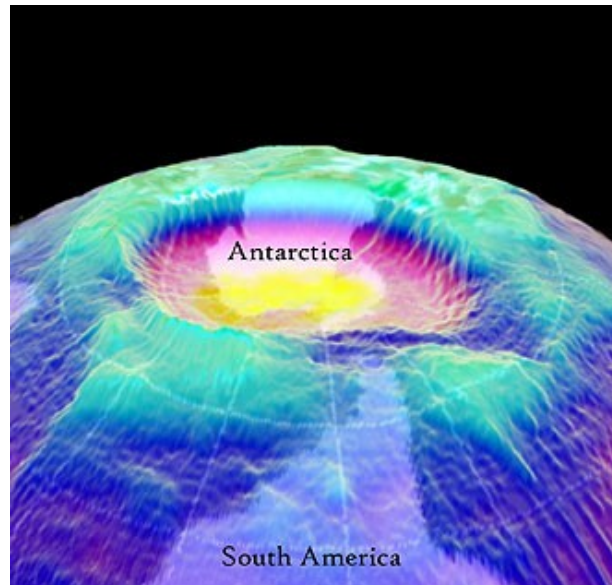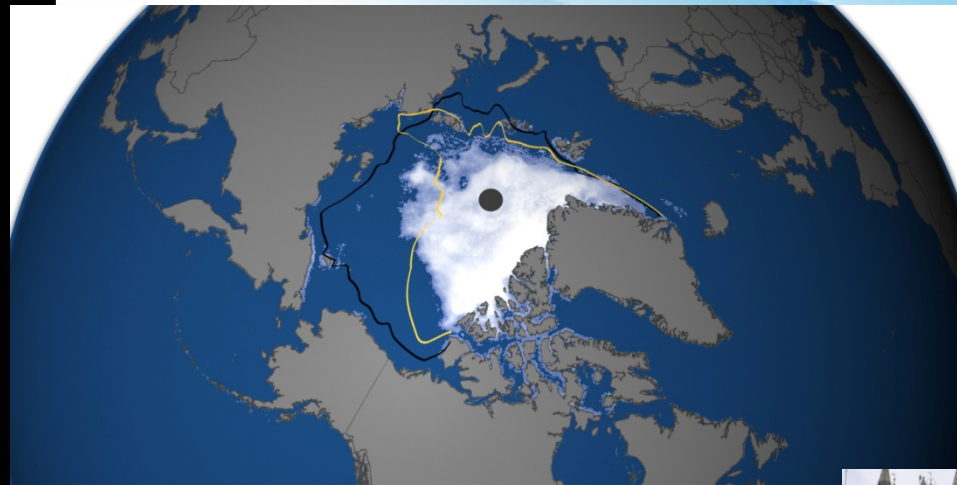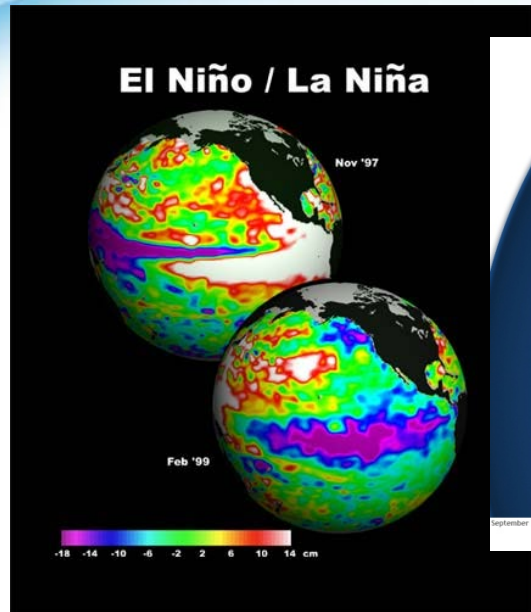Mohan Ramamurthy
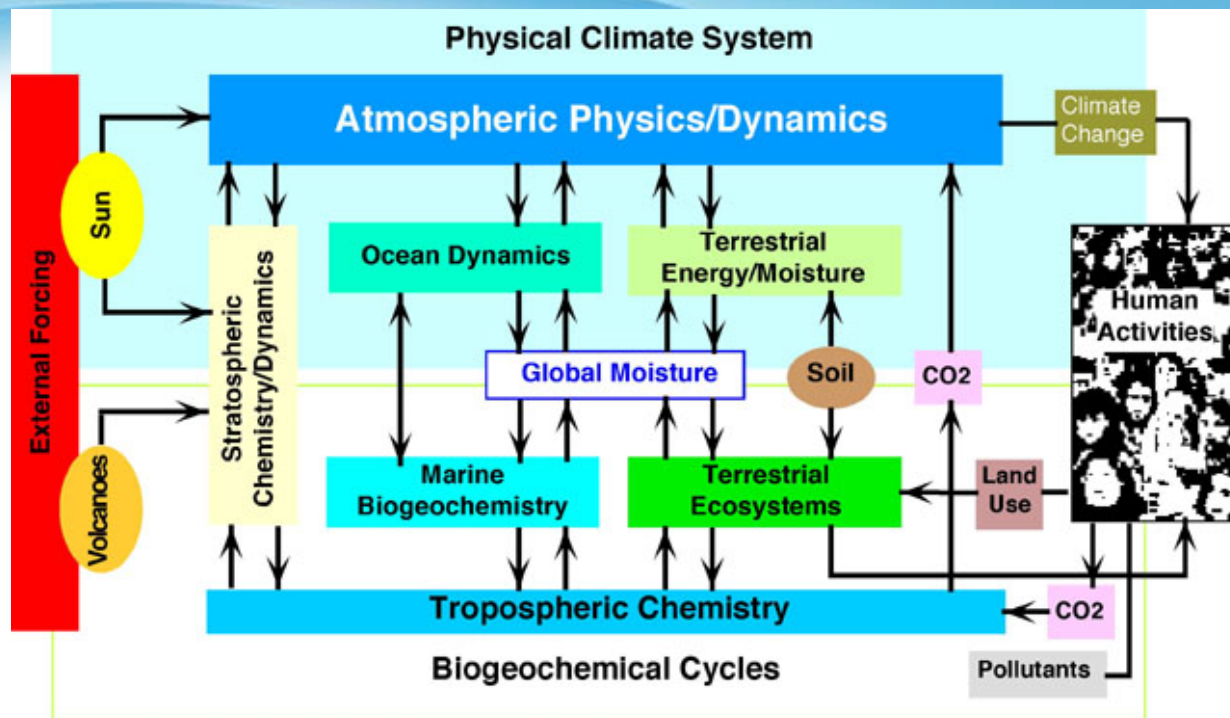Unidata Program Center
UCAR Community Programs
Boulder, CO

# Science and Society



It is estimated that at least 1/3 of the U.S. GDP is weather & climate sensitive, a potential impact of $4 trillion/year.

CLIMATE CHANGE

# Grand Challenges are Global & Multidisciplinary

# Earth System Science



(from Earth System Science: An Overview, NASA, 1988)

Need Cyberinfrastructure to support ESS thinking.

It requires data and information *integration* and knowledge *synthesis* across "systems" or domains.

Challenge: Providing the right data, in the right format, to the right application.

# Networked Science

➢ Distributed knowledge communities working collaboratively as a virtual community to tackle problems never before possible.
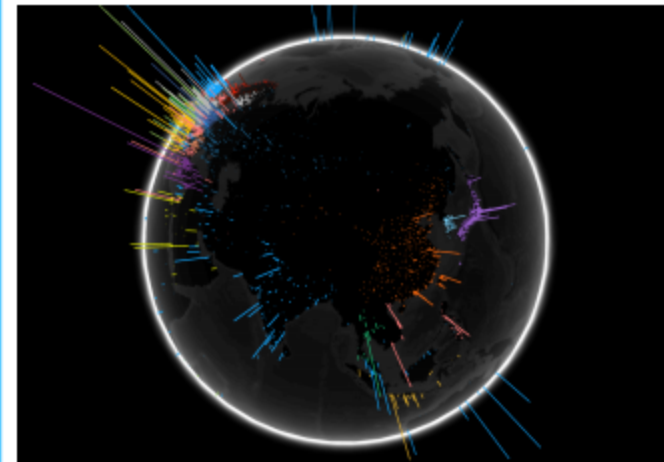
# BIG DATA BUZZ

# Digital Universe

According to a study by IDC, 1.8 Zettabytes of information was created in 2011.



From Science Projects to Social Networks to Smart Technology

MORE DEVICES
MORE APPLICATIONS
MORE CONTENT
MORE ON-DEMAND ACCESS

Finding Answers where there are yet to be Questions...

Source: IDC's Digital Universe Study, sponsored by EMC, June 2011



7910 EXABYTES

A DECADE OF DIGITAL UNIVERSE GROWTH

1227 EXABYTES

130 EXABYTES

2005    2010    2015

# The Era of Data-Intensive Science

Data is the lifeblood of science, but we need to move from creating data to discovering knowledge.

# "Sea of Data"

GOES-R (2017)

JPSS (2014)

~3 Tb of data/day

Phased Array Radar, with 20 to 30-second volume scans, compared with 5-7 mins. with current radars.

Global, high-resolution coupled models integrated in ensemble mode from days to decades

# CMIP-5 & IPCC Fifth Assessment



## Introduction to CMIP5:
## The Experiments

Example: CMIP5 long-term suite of experiments



Green subset is for coupled carbon-cycle climate models only

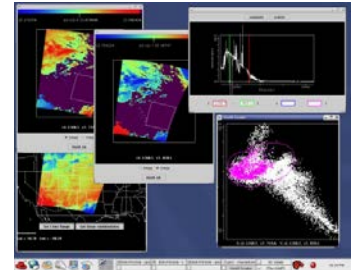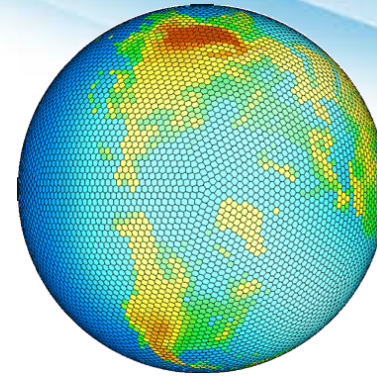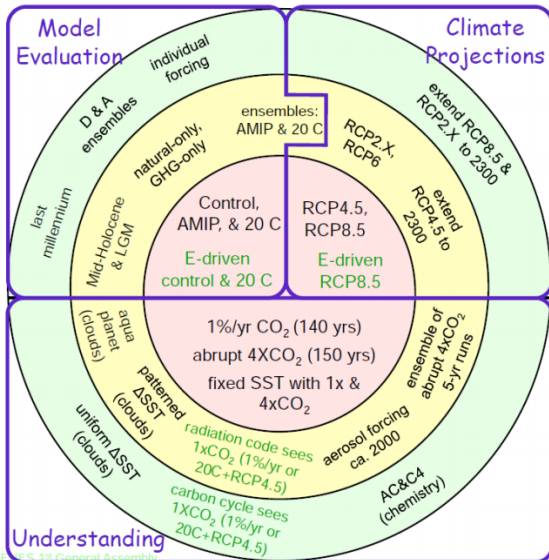| Core: | ≥1718 yrs |
| Tier 1: | ≥1727 yrs |
| Tier 2: | ≥2038 yrs |

K. E. Taylor

## CMIP5 in numbers

**Simulations:**

~90,000 years

~60 experiments

~20 modelling centres using

~30 major(*) model configurations

~2 million output datasets

~10's of petabytes of output

~2 petabytes of CMIP5 requested output

~1 petabyte of CMIP5 "replicated" output

~10 TB of land-biochemistry (from the long term experiments alone).

**Of the replicants:**

~ 220 TB decadal

~ 540 TB long term

~ 220 TB atmos-only

~100 TB of 3hourly atmos data!

~215 TB of ocean 3d monthly data!

~250 TB for the cloud feedbacks!

**Expected Usage (@ BADC):**

~ hundreds of users downloading at a sustained daily average rate of between 1 and 3 Gbit/s (or up to 35 TB/day from BADC ...)

Source: Bryan Lawrence, British Atmospheric Data Centre

# Expected Increase in Data Volume



Source: NCDC, NOAA

# A Provocative Suggestion



**WIRED**

The End of Theory: The Data Deluge Makes the Scientific Method Obsolete

By Chris Anderson ✉    06.23.08

Illustration: Marian Bantjes

Wired, 23 June 2008 issue

# Data Challenges: The Five "V"s

- **Volume:** Explosion of data
- **Variety:** Different types of data (e.g, multidisciplinary, societal information, etc.); interoperability
- **Velocity:** Speed of discovery, access, analysis, integration, and visualization;
- **Views:** Many consumers of data (e.g., researchers, educators, students, policy makers, social scientists, and the public); Diverse applications; Multiple devices;
- **Virtual Communities:** Globally distributed, different cultures, practices, and policies

# The Long Tail Problem



National Science Foundation 2007 Awards

Heidorn, P. Bryan (2008). Shedding Light on the Dark Data in the Long Tail of Science. Library Trends 57(2) Fall 2008 .

# The Long Tail Problem - contd

- By some estimates, only 5% of the data generated by individual PIs is shared
- There are many reasons for this – both sociological as well as technical
  - Lack of incentives
  - Inadequate resources
  - Additional burden or unfunded mandate
  - Protectiveness – PIs don't want to be "scooped"
  - Technical challenges
  - Absence of local or community data repositories
- Need to give PIs the tools required for sharing their data; also need tools for adding metadata
  - Need to create incentives
  - Need to change the culture

# NSF Data Management Plan

- <u>All proposals must include </u>a Data Management Plan.

- Plan should describe how the proposal will conform to NSF policy on dissemination and sharing of research results.

- Plan will be reviewed as part of the intellectual merit and/or broader impacts of the proposal depending on the proposal intent.

- NSF will not permit submission of a proposal that is missing a data management plan.

# Data Citation: The Next Frontier?

- Scientific publications should be accompanied by data, algorithms, models, and parameters – need comprehensive data citation. Need transparency. Important for reproducibility.

- This is not just a technical challenge, but it is also a major cultural and organizational challenge.
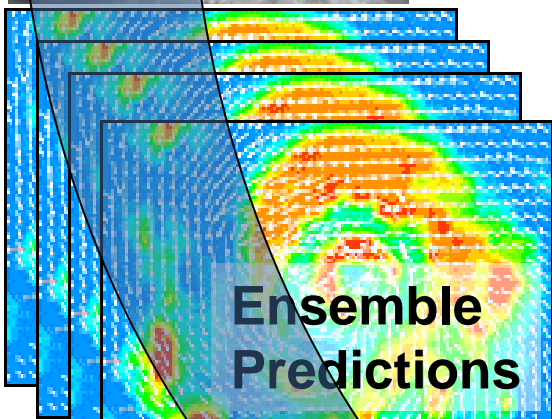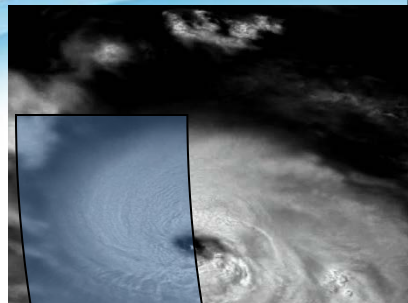
## Proposed AMS Statement on the Importance of Data Availability in Support of Scholarly Publications

Repeatability of research results is one of the main ten...
science; it is a benchmark upon which the reliabili...
placed in authors' conclusions. Some d...
generated them, and making ...
will further scientif...
independ...
...own false
...cy for its scholarly

■ ...ors be required to agree to make data and
...promptly available to readers upon request, subject to
...aints. Exceptions may be appropriate in certain limited
...stances to preserve privacy, to assure patent protection, or for other legal
reasons. Any restrictions on the availability of materials or information must be
disclosed to the Editor at the time of manuscript submission.

2. Data sets must be made available to Editors and peer-reviewers if required at the
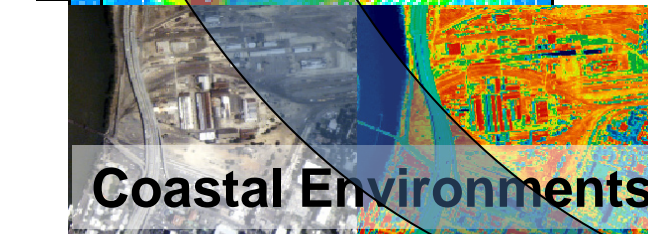time of review in order to ensure a comprehensive peer-review process. The AMS

"We have never done it that way before," should be academia's motto, said Kathleen Fitzpatrick, a professor of media studies at Pomona College.

# End-to-End Data Services Chained by Workflows

Understanding societal impact of flooding from hurricanes involves integrating data from atmospheric sciences, oceanography, hydrology, geology, and social sciences and interfacing the results with decision support systems.

# GIS Integration: An Enormous Opportunity

- Need geospatially-enabled cyberinfrastructure so that information can be integrated for location-based understanding of events, processes, interactions, and impacts. (May Yuan)

- GIS integration should not be an after thought. Scientific data systems need to directly enable GIS tools.

# GIS and THREDDS Data Server working together

UCAR

Unidata

**NCAR model air temperature anomaly**

- Data from NCAR climate model
- TDS software from Unidata
- WMS software from University of Reading
- Server run by NCAR GIS program
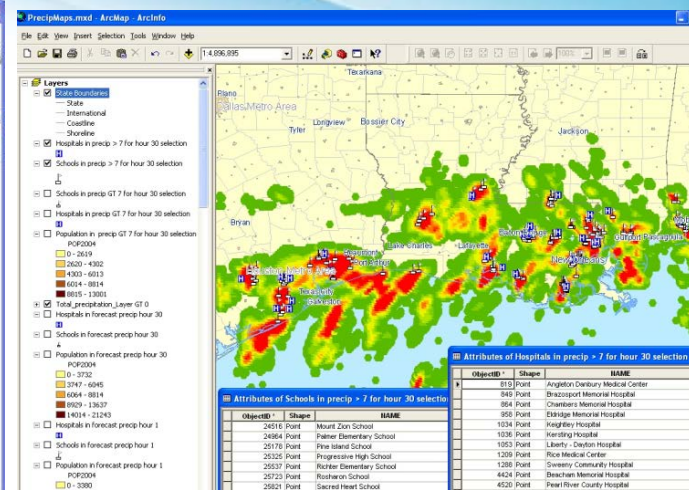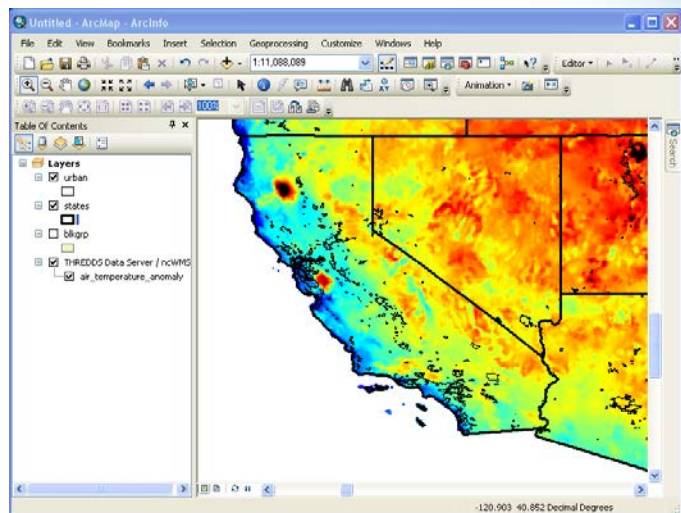- Map and analysis software from ESRI

**NAM output displayed by ArcGIS**



**WRF output displayed by Google Maps**

# "Cloud" Computing

- Data volumes are too large to bring all of the data to your local environment

- Need to keep them close to the point of origin or dissemination

- Will need to move more of the processing, applications, and computations on to the server (e.g., GDS, NCO, etc.).

- Impractical to store/serve data in multiple formats – so need built-in translators, brokers and mediation services.

- Industry is well ahead of the science world in dealing with Big Data, Cloud Computing, and Virtualization (e.g., SaaS, Application Service Providers, MapReduce, Hadoop, etc.)

# Virtualization

- **Virtualization** is a tool that has many technical uses, most of which have nothing to do with the cloud.

- Virtualization allows the use of a single piece of physical hardware, to perform the work of many. e.g., multiple operating system instances running on one hardware device are far more economical than a single piece of hardware for every server task.

- You can do virtualization on your laptop or desktop, although it is unlikely that either will be a cloud server.

# Mobile Computing Devices

- Smartphones, tablets, laptops, etc.

- Ubiquitous access to network (anywhere, any place, etc.), with occasional offline needs

- Tablet applications, using thin client-server approaches, are currently being developed by ESRI, NWS, and others.
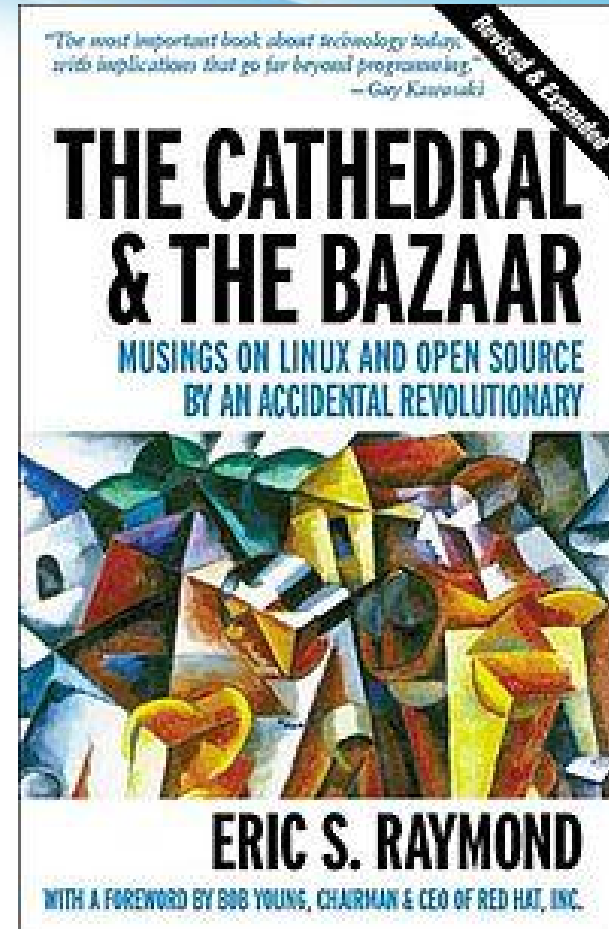


## AWIPS II Thin Client

- Primary users: Incident METeorologists
  - NWS forecasters on site at fires and other emergencies
- Laptop computers powerful enough already
  - PADS project
- Web services reduce data requirements
  - Specify areal subsets
  - Minimizes comms requirements

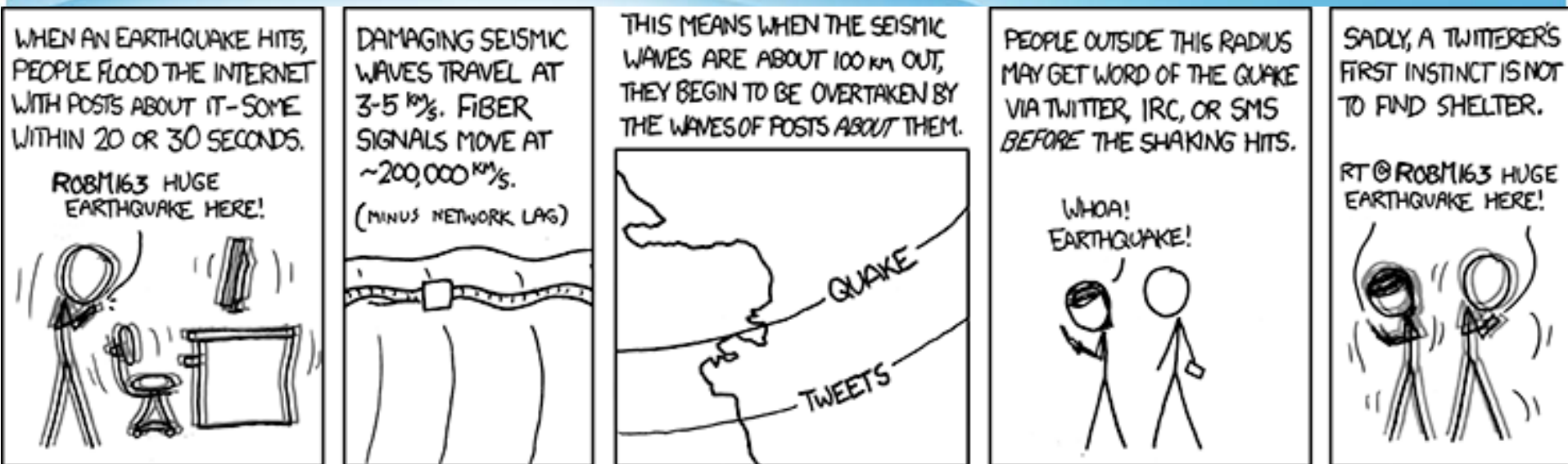# Open Source Software Development in the Geosciences

- The software community has many successful examples of Open Source Software Development (Linux, Mozilla, Apache,

- The atmospheric science community has had many decades of experience with community modeling efforts (NCAR GCM, CCSM, CESM, MM5, and WRF).

- Beyond that, there has been little attention paid to Open Source Development in our field.
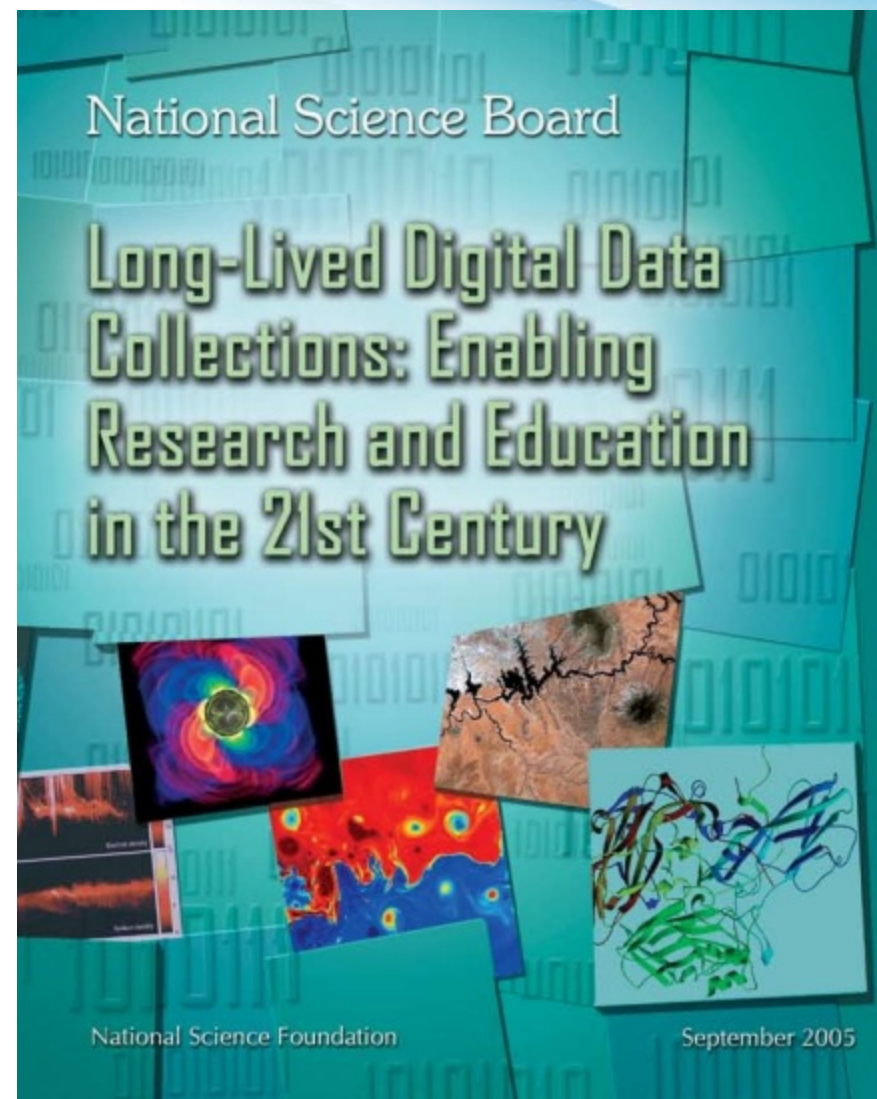
**Wikis, Blogs, Twitter, Facebook, and Linked-In in increasing use**

**Mobile sensors and crowd sourcing and validation, using social network**

**Beyond being a communication tool, how do you use Social Networking for data collection, aggregation, and sharing?**

# Education, Training and Workforce Development

- Education and workforce requirements are not aligned;

- Computational, geospatially, and data literacy needs to be advanced;

- Far too few women and minorities in IT disciplines and in the workforce;

National Science Board

Long-Lived Digital Data Collections: Enabling Research and Education in the 21st Century

National Science Foundation

September 2005

# Unidata Strategic Plan

**Unidata 2020:** **Geoscience at the Speed of Thought** *through accelerated data discovery, access, analysis, and visualization.*

**Mission:** **To transform the geosciences community, research, and education by providing innovative data services and tools.**

Reduce *"data friction",* lower the barriers, and reduce "time to research*"*

Accelerate *user workflows (manual or automated)*

Contribute toward flipping the 4:1 ratio (the current 80-20 situation)

# Concluding Remarks

❖ We live in an exciting era in which advances in computing and communication technologies, coupled with a new generation of geoinformatics, are accelerating scientific research, creating new knowledge, and leading to new discoveries at an unprecedented rate.

❖ Workshops like this play an ever more important role in bringing people together to understand our community needs, examine progress, look at opportunities, address challenges, and foster new partnerships.

❖ Thank you!