# Building an Ecosystem to Accelerate Data-Driven Innovation
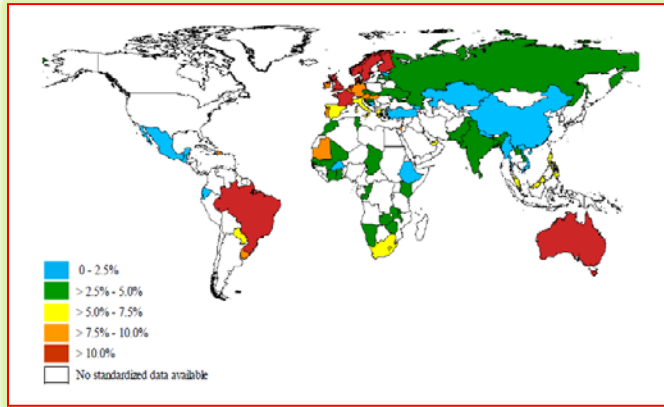
**Dr. Francine Berman**

Chair, Research Data Alliance / US

Edward P. Hamilton Distinguished Professor of Computer Science, Rensselaer Polytechnic Institute
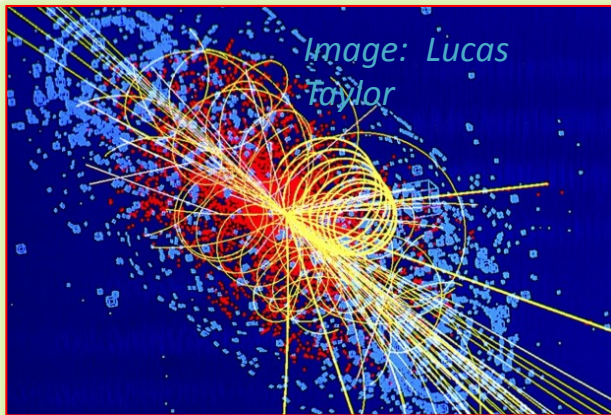
Fran Berman

# Research Data Driving Solutions to Complex Scientific and Societal Challenges



**Who is most at risk to contract asthma?**



**How can we increase wheat yields?**



Image: Lucas Taylor

**How accurate is the Standard Model of Physics?**



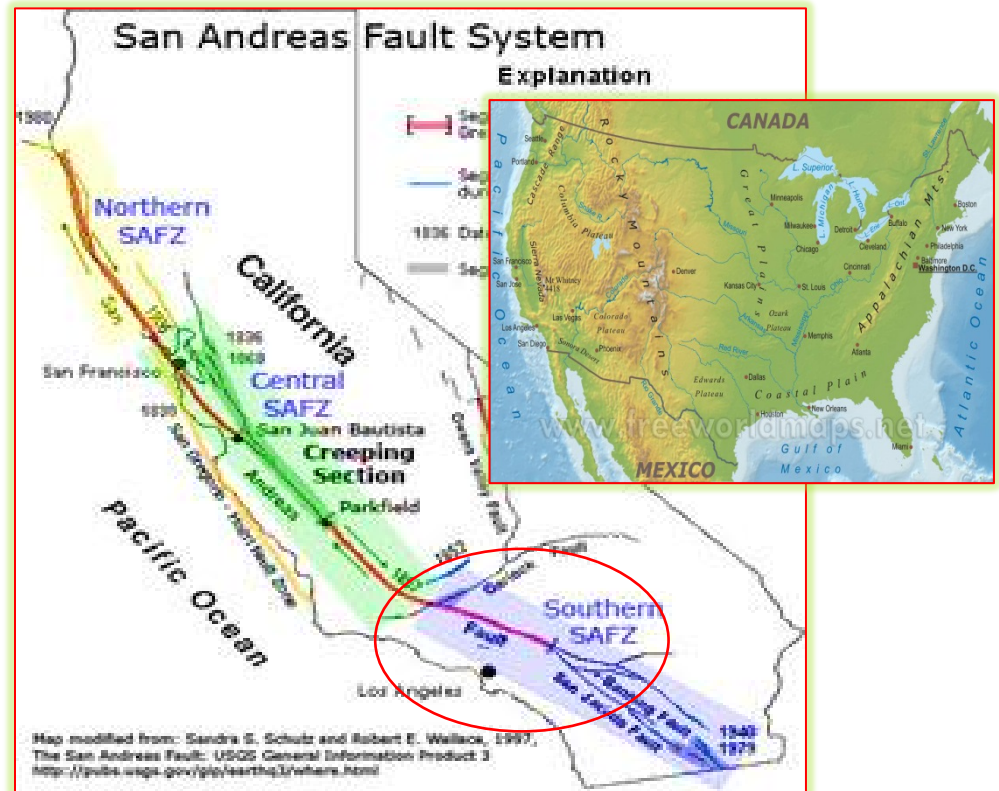Image: Ceinturion, Wikipedia

**How can we best address energy needs and sustain the environment?**

# Data-Driven Geoscience: How Can We Respond to Large-Scale Earthquakes?

## Earthquake simulations enable

- Enhanced **scientific understanding** of the physical world

- More strategic plans for bridge, building and other physical infrastructure reinforcements to **increase safety**

- Better **disaster response planning** for police, fire fighters, ER teams in high-risk areas to increase their effectiveness



*Simulation courtesy of Amit Chourasia, SDSC, Table information courtesy of Southern California Earthquake Center*

Fran Berman

# TeraShake Simulation of 7.7 Earthquake on the Lower San Andreas Fault

## Earthquake simulations enable

- Enhanced **scientific understanding** of the physical world

- More strategic plans for bridge, building and other physical infrastructure reinforcements to **increase safety**

- Better **disaster response planning** for police, fire fighters, ER teams in high-risk areas to increase their effectiveness
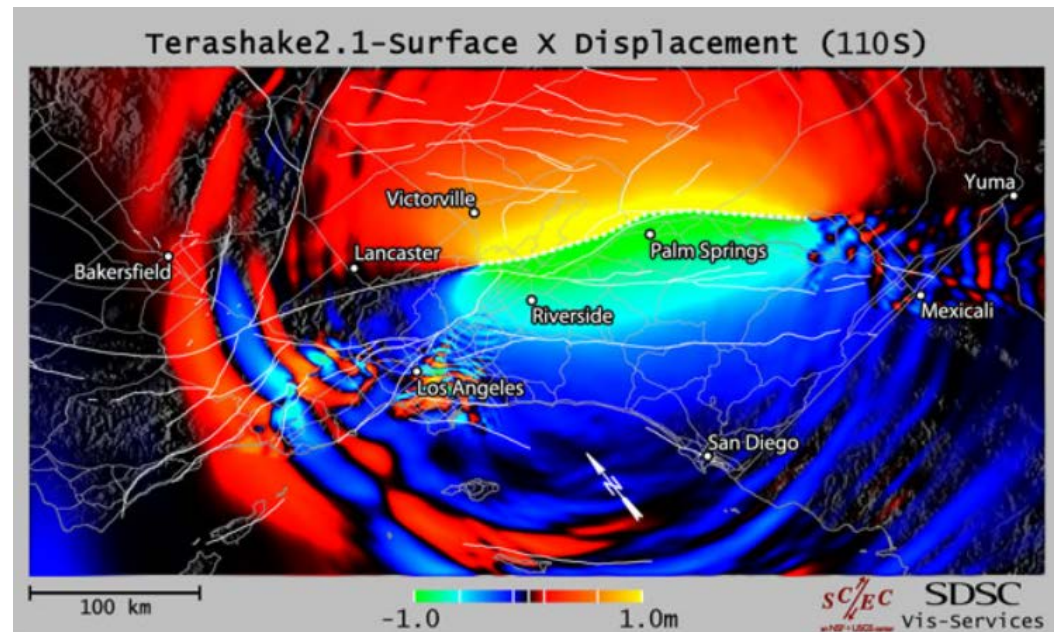


*Simulation courtesy of Amit Chourasia, SDSC, Table information courtesy of Southern California Earthquake Center*

Fran Berman

# More Accurate Simulations Require More Infrastructure

## Earthquake simulations enable

- Enhanced **scientific understanding** of the physical world

- More strategic plans for bridge, building and other physical infrastructure reinforcements to **increase safety**

- Better **disaster response planning** for police, fire fighters, ER teams in high-risk areas to increase their effectiveness

| Estimated figures for simulated 240 second period, 100 hour run-time | **TERASCALE:** TeraShake domain (600x300x80 km^3) | **PETASCALE:** PetaShake domain (800x400x100 km^3) |
|---|---|---|
| **Fault system interaction** | NO | YES |
| **Inner Scale** | 200m | 25m |
| **Resolution of terrain grid** | 1.8 billion mesh points | 2.0 trillion mesh points |
| **Magnitude of Earthquake** | 7.7 | 8.1 |
| **Time steps** | 20,000 (.012 sec/step) | 160,000 (.0015 sec/step) |
| **Surface data** | **1.1 TB** | **1.2 PB** |
| **Volume data** | **43 TB** | **4.9 PB** |

*Simulation courtesy of Amit Chourasia, SDSC, Table information courtesy of* Fran Berman
*Southern California Earthquake Center*

# Integrated Infrastructure Critical:
## Application Needs Span the Spectrum

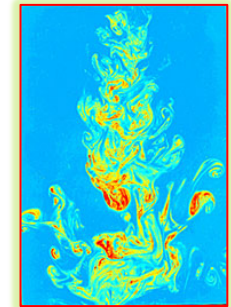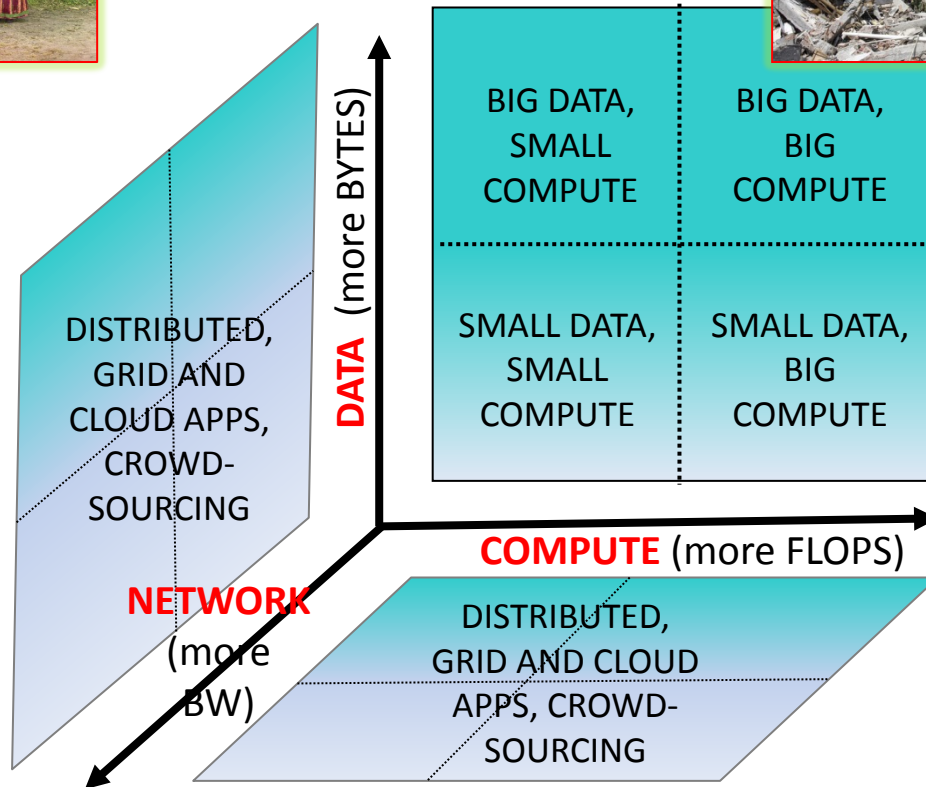*Digital history image integration*

*Earthquake Simulation*

*Cosmology*

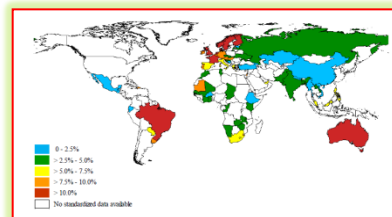*Analysis and modeling of protein function and structures*

DISTRIBUTED, GRID AND CLOUD APPS, CROWD-SOURCING

**DATA** (more BYTES)

| BIG DATA, SMALL COMPUTE | BIG DATA, BIG COMPUTE |
| SMALL DATA, SMALL COMPUTE | SMALL DATA, BIG COMPUTE |

**COMPUTE** (more FLOPS)

**NETWORK** (more BW)

DISTRIBUTED, GRID AND CLOUD APPS, CROWD-SOURCING

*Turbulence*

*Analysis of Data from the CERN Large Hadron Collider*
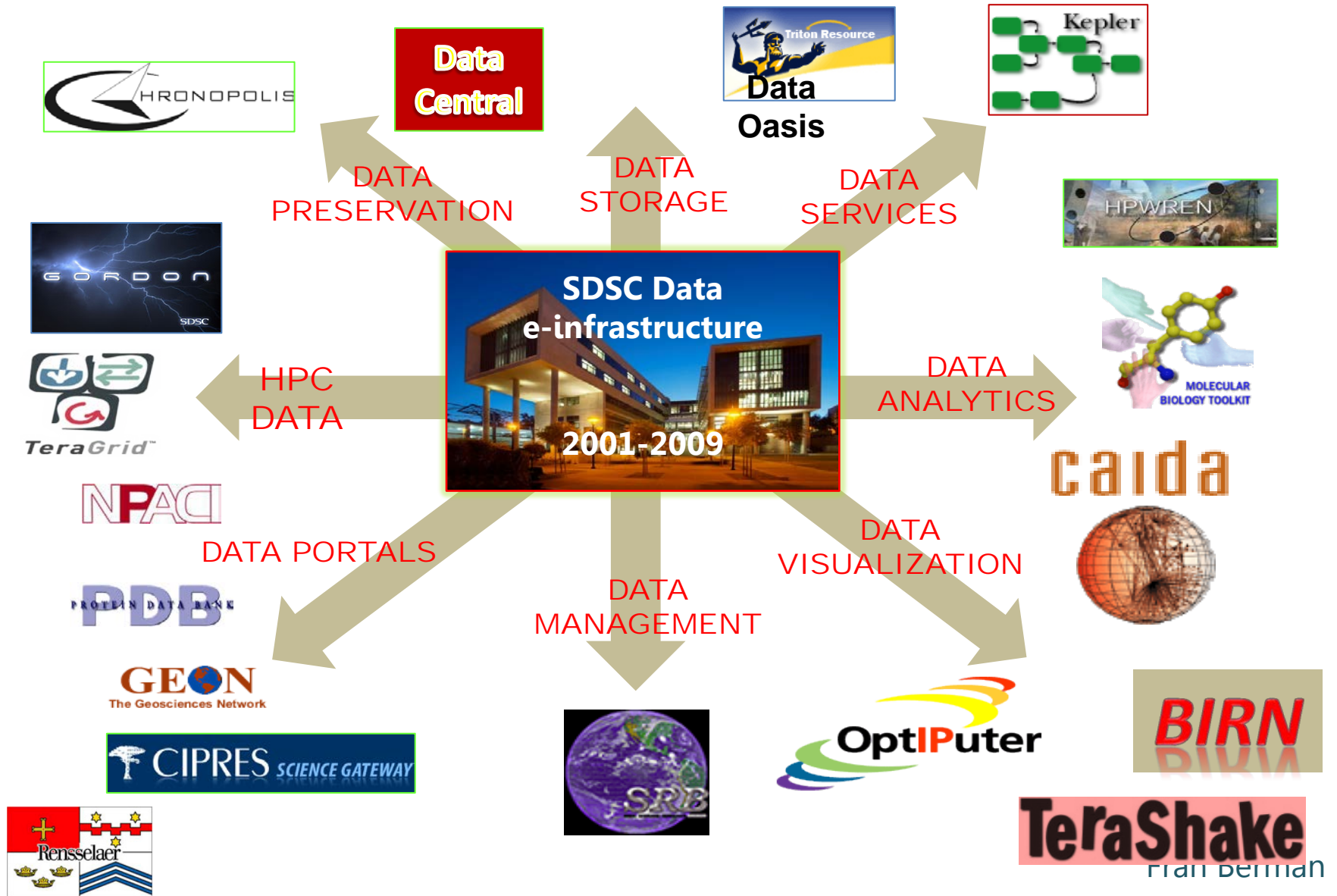
*Global Public health*

*Seti@home, MilkyWay@Home, BOINC*

# Many Kinds of Technical Data Infrastructure Needed to Drive Innovation

# Social, Organizational, and Human Infrastructure Equally Important

**Social and Organizational Infrastructure**

**Human Infrastructure / Workforce**



€250 billion
potential annual value to Europe's public sector administration—more than GDP of Greece

$600 billion
potential annual consumer surplus from using personal location data globally

60% potential increase in retailers' operating margins possible with big data

140,000–190,000
more deep analytical talent positions, and

1.5 million
more data-savvy managers needed to take full advantage of big data in the United States
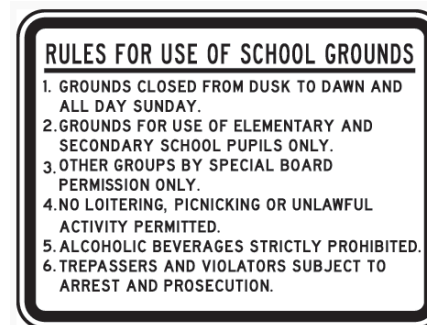
*Data Scientists*

*Data-focused Curriculum and Training*

*Community Practice*

*Sustainable Economics*

RULES FOR USE OF SCHOOL GROUNDS
1. GROUNDS CLOSED FROM DUSK TO DAWN AND ALL DAY SUNDAY.
2. GROUNDS FOR USE OF ELEMENTARY AND SECONDARY SCHOOL PUPILS ONLY.
3. OTHER GROUPS BY SPECIAL BOARD PERMISSION ONLY.
4. NO LOITERING, PICNICKING OR UNLAWFUL ACTIVITY PERMITTED.
5. ALCOHOLIC BEVERAGES STRICTLY PROHIBITED.
6. TRESPASSERS AND VIOLATORS SUBJECT TO ARREST AND PROSECUTION.

*Policy*

*Common Standards*

# Today's Presentation:
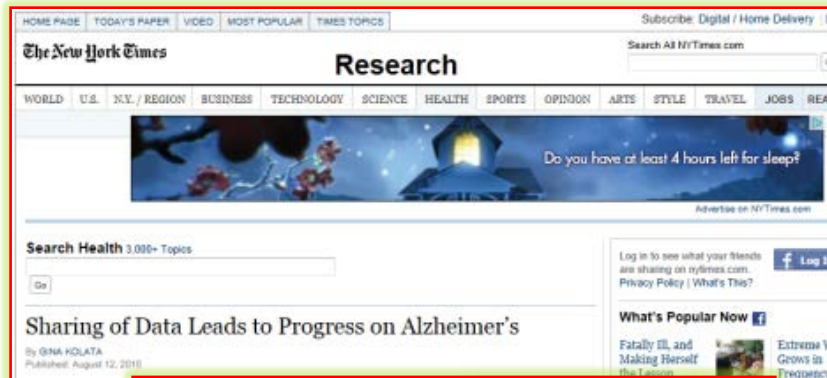# Two perspectives on Research Data Infrastructure

- **Opportunity:** **Maximizing data-driven innovation through data sharing and exchange**

  – Efforts to build a coordinated global Infrastructure to support data access, sharing and use

- **Challenge:** **Prioritizing the Development, Implementation and Sustainable Support of Data Infrastructure**

  – Strategies to accelerate efforts within organizations, communities and sectors

Fran Berman

# Opportunity:  Maximizing Data-Driven Innovation Through Data Sharing and Exchange



Fran Berman

# Data-Sharing Driving innovation Across Sectors and Communities



The New York Times — Research

**Sharing of Data Leads to Progress on Alzheimer's**
By GINA KOLATA
Published: August 12, 2010

**iHealthBeat** — Reporting Technology's Impact on Health Care

NEWS ARCHIVE

**Report: Data Sharing Pilot Helped Hospitals Save $11B, 136K Lives**
Tuesday, February 18, 2014

OFFICE OF JUSTICE PROGRAMS

**NATIONAL INSTITUTE OF JUSTICE**
Research • Development • Evaluation

NIJ JOURNAL NO. 267

**In Brief: Expanding Research by Sharing Data**
by NIJ staff
*NIJ makes data available for future research.*

InformationWeek **Healthcare**

NEWS
**Sharing Psychiatry EHR Data Cuts Readmission Rates**

**theguardian**

News | US | World | Sports | Comment | Culture | Business | Money
Professional › Global Development Professionals Network › Sign

**Farming and food security hub**
From the Global Development Professionals Network

**How might open data in agriculture help achieve food security?**
The policy support for improving the ability to store and share data on agriculture is growing. But how do you ensure farmers in developing countries benefit and will it achieve food security?

Caspar van Vark
Guardian Professional, Monday 25 November 2013 06.09 EST
Jump to comments (1)

**Astronomers Release Unprecedented Data Set on Celestial Objects that Brighten and Dim**

PASADENA, Calif.—Astronomers from the California Institute of Technology (Caltech) and the University of Arizona have released the largest data set ever

# World-wide Efforts Focusing on Infrastructure to Support Research Data Sharing, Access, Use



A Europe-Japan-United States GNSS data-sharing pilot project for the Geohazard Supersites and Natural Laboratories

Falk Amelung, University of Miami, USA (GEO task lead)
Craig Dobson, NASA and Committee of Earth Observation Satellites (CEOS)
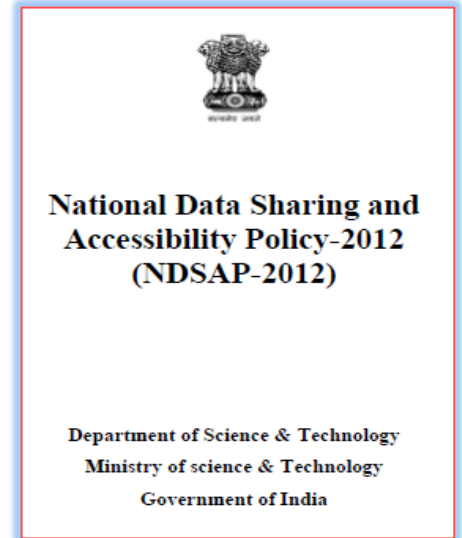Rui Fernandes, EPOS and EUREF <rmanuel@di.ubi.pt>



SCIENCE BUSINESS

Policy Analysis: Investment, Policy

SPECIAL COVERAGE HORIZON 2020

Promote data sharing to advance global research say policy leaders

Richard L. Hudson, Science|Business

EU and US experts see big benefits from scientists sharing more data - but say global agreement on privacy, literacy and other issues is needed

**Science, Humanities, Arts Communities**

**Libraries, Archives, Repositories, Museums**



DATA.GOV
EMPOWERING PEOPLE

**E-Infrastructure professionals, data analysts, data center staff, …**



National Data Sharing and Accessibility Policy-2012 (NDSAP-2012)

Department of Science & Technology
Ministry of science & Technology
Government of India

**Data Scientists**

Australian National Data Service

Our Vision: More Australian researchers reusing research data more often

ANDS is enabling the transformation of:

Data that are:    to    Structured Collections that are:
Unmanaged    →    Managed
Disconnected    →    Connected
Invisible    →    Findable

ANDS News

ANDS webinar
Join ANDS to discuss national and international trends in research data infrastructure

Congratulations
The Public Record Office of Victoria has completed its ANDS-funded project

eResearch Australasia
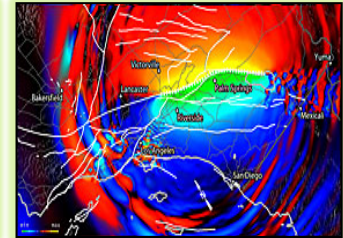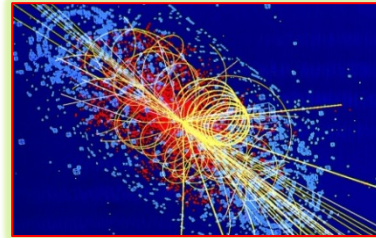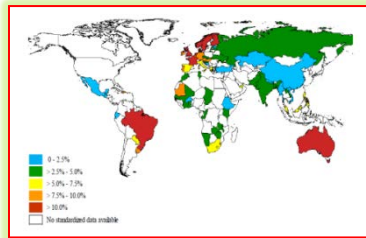
Fran Berman

# Many Infrastructure Building Blocks Needed to Accelerate Progress



| Research Dissemination and Reproducibility | Data Use and Re-use | Data Access (now) and Preservation (later) | Data Discovery and Data Sharing |
| --- | --- | --- | --- |

| Data Access and Distribution Policy | Institutional Data Sharing Practice | Data Discovery Tools |
| --- | --- | --- |
| Digital Object Identifiers | Common Metadata Standards | Data Citation Standards |
| Data Preservation Practice | Data Analytics Algorithms | Data Scientists and Expert Support |
| Curation Practice and Policy | Sustainable Economic Models | Auditing, Certification and Reporting Practice |

Fran Berman

# Research Data Alliance Created to Accelerate Development of Research Data Sharing Infrastructure Worldwide

- RDA is an emerging, global community-driven organization created to **accelerate the development of research data-sharing infrastructure** world-wide.

- RDA community efforts focus on building **social, organizational and technical infrastructure** to

  - *reduce barriers to data sharing and exchange*

  - *accelerate the development of coordinated global data infrastructure*





Fran Berman

# RDA Approach:
## CREATE → ADOPT → USE



**RDA Members come together as**

- **Working Groups** – 12-18 month efforts to build, adopt, and use specific pieces of infrastructure

- **Interest Groups** – longer-lived discussion forums that spawn Working Groups as specific pieces of needed infrastructure are identified.

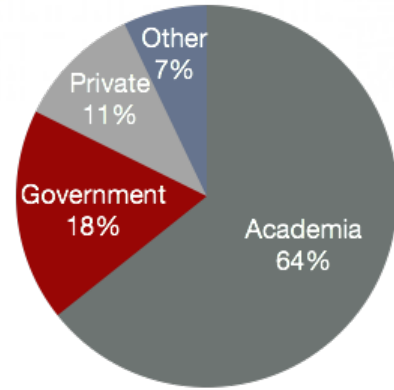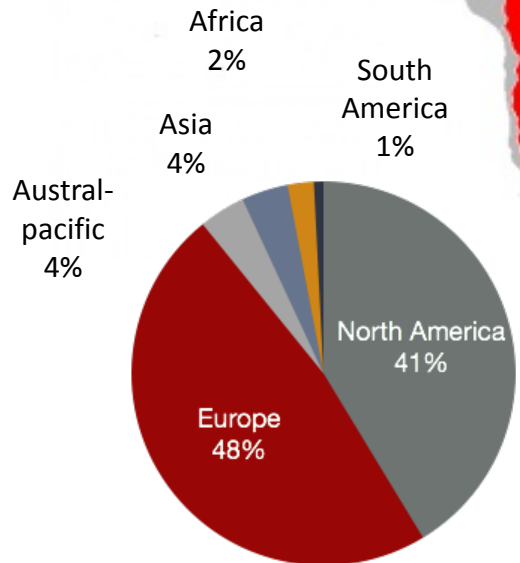**Working Group efforts focus on the development and use of data sharing infrastructure**

- **Code, policy, infrastructure, standards, or best practices that are adopted and used** by communities to enable data sharing

- **"Harvestable" efforts** for which 12-18 months of work can eliminate a roadblock

- **Efforts that have substantive applicability** to groups within the data community, but may not apply to everyone

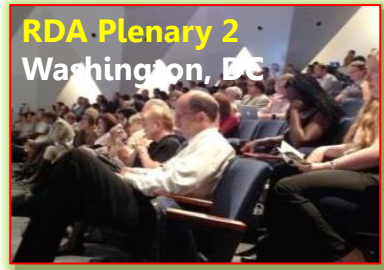- **Efforts for which working scientists and researchers can start today**

Fran Berman

# The RDA Community Today: Over 1600 members from 70+ countries (as of 15/3/14)



Map courtesy traveltip.org

**Geographic distribution:**
- North America 41%
- Europe 48%
- Austral-pacific 4%
- Asia 4%
- Africa 2%
- South America 1%

**Sector distribution:**
- Academia 64%
- Government 18%
- Private 11%
- Other 7%

# Precipitous Growth


**RDA Plenary 1 / Launch**
Gothenburg, Sweden


**RDA Plenary 2**
Washington, DC


**RDA Plenary 3**
Dublin, Ireland


**Amsterdam**

**First "neutral space" community meeting (Data Citation Summit)**

**First Org. Partner Meet-up**

**First BOFs**

**380 participants from 22 countries**

**First Organizational Assembly**

**6 co-located events**

**14 BOF, 12 Working Groups, 22 Interest Groups**

**497 participants**

**First Working Groups and Interest Groups**

**240 participants**

Global Data Planning Meeting: October 2012

**RDA Launch / First Plenary**

**RDA Second Plenary**

**RDA Third Plenary**

**RDA Fourth Plenary**

**March 2013**    **September 2013**    **March 2014**    **September 2014**

First RDA organizational telecon: August 2012

First Working Group exchange meeting

Fran Berman

# RDA Plenary 3
## March 26-28, 2014
## Dublin, Ireland



jenniferlin15
@jenniferlin15

Follow

Ferguson: intn'l collaboration as critical criteria for evaluating research impact. #altmetrics provide otherwise invisible view #RDAPlenary

5:50 AM - 26 Mar 2014

2 RETWEETS

**6 co-located data-focused events**

**Organizational Assembly Meeting: 20+ new members**

| Professional Title | Total | % |
|---|---|---|
| Advisor/Consultant | 22 | 4% |
| CEO / Managing Director / Chief Executive | 35 | 7% |
| CTO / IT Director | 20 | 4% |
| IT Specialist / IT Architect | 53 | 11% |
| Journalist / Editor / Copywriter | 6 | 1% |
| Librarian | 27 | 5% |
| Other | 93 | 19% |
| Student | 38 | 8% |
| Policy Development Manager / Policy Consultant | 12 | 2% |
| Professor | 42 | 8% |
| Programme Manager / Project Manager | 62 | 12% |
| Researcher | 87 | 18% |
| Total | 497 | 100% |

Fran Berman

# RDA Interest (IG) and Working Groups (WG) by Focus
## (as of 15/3/14)

## Domain Science - focused

- Toxicogenomics Interoperability IG
- Structural Biology IG
- Biodiversity Data Integration IG
- Agricultural Data
- Interoperability IG
- Digital History and Ethnography IG
- Defining Urban Data Exchange for Science IG
- Marine Data Harmonization IG
- Materials Data Management IG

## Community Needs - focused

- Community Capability Model IG
- Engagement IG
- Clouds in Developing Countries IG

## Reference and Sharing - focused

- Data Citation IG
- Data Categories and Codes WG
- Legal Interoperability IG

## Data Stewardship - focused

- Research Data Provenance IG
- Certification of Digital Repositories IG
- Preservation e-infrastructure
- Long-tail of Research Data IG
- Publishing Data IG
- Domain Repositories IG
- Global Registry of Trusted Data Repositories and Services IG

## Base Infrastructure - focused

- Data Foundations and Terminology WG
- Metadata Standards WG
- Practical Policy WG
- PID Information Types WG
- Data Type Registries WG
- Metadata IG
- Big Data Analytics IG
- Data Brokering IG

**HPC Members welcome!**

# First RDA Infrastructure Deliverables coming this Fall

## Data Type Registries WG

- **Deliverables:** System of data type registries, formal model for describing types, working model of a registry.
- **Initial Adopters and Users:** CNRI, International DOI Foundation, Deep Carbon Observatory

## Practical Code Policies

- **Deliverables:** Survey of policies in production use, testbed of machine actionable policies, deployment of 5 policy sets, policy starter kits
- **Initial Adopters and Users:** RENCI, DataNet Federation Consortium, CESNET, Odum Institute, EUDAT

## Persistent Identifier Information Types

- **Deliverables:** Minimal set of PID types, API
- **Initial Adopters and Users:** Data Conservancy, DKRZ

## Language Codes

- **Deliverables:** Operationalization of ISO language categories for repositories.
- **Initial Adopters and Users:** Language Archive, Paradisec

## Data Foundations and Terminology

- **Deliverables:** Common vocabulary for data terms, formal definitions and open registry for data terms
- **Initial Adopters and Users:** EUDAT, DKRZ, Deep Carbon Observatory, CLARIN, EPOS

## Metadata Standards

- **Deliverables:** Use cases and prototype directory of current metadata standards starting from DCC directory
- **Initial Adopters and Users:** JISC, DataOne

# Next Steps for the RDA

**More Infrastructure** → **Continuing pipeline of infrastructure deliverables adopted and used to accelerate data sharing**

**Increasing coordination of infrastructure**

**Effective Community** → **Increasing cross-boundary collaborations between domains, sectors, organizations**

**Synergistic Programs** → **International and regional programs focusing on workforce, outreach, expansion of infrastructure impact**

**Focus on Industry** → **New partners in the Organizational Assembly**

**Focused strategy to support development of industry infrastructure for data sharing**

Fran Berman

# Challenge: Supporting and Sustaining Research Data Infrastructure



Fran Berman

# Research Data and Data Sharing Key to Innovation.
## Research Data an Accelerator for All Sectors.

- National governments increasingly calling for public access to research data.

  - *Valuable to all sectors as a driver of current and future innovation*

- **Research data infrastructure is necessary** to support

  - *Use and re-use*

  - *Research reproducibility*

  - *Federally mandated data management plans*

  - *Public access to research data*



DIGITAL AGENDA FOR EUROPE
A Europe 2020 Initiative

Digital Agenda for Europe



Sharing of Data Leads to Progress on Alzheimer's



Report: Data Sharing Pilot Helped Hospitals Save $11B, 136K Lives



Farming and food security hub
From the Global Development Professionals Network

How might open data in agriculture help achieve food security?

Fran Berman

# Yet Research Data Infrastructure is Difficult to Sustain. Why?



Data-at-Risk & Rescue Inventory

:: Browse Items    :: Browse Collections

Documenting Scientific Data-at-Risk and Data Rescue

**What is the Data-at-Risk Inventory (DARI)?**

The Data-at-Risk & Rescue Inventory (DARI) is an initiative that:
1. Catalogs valuable scientific data that are at risk of being lost to posterity.

---

The Sydney Morning Herald

**Federal Politics**

Some apps sold separately; vary by market.

**Political News**  Political Opinion  **Breaking Politics Video**  The Pulse  The Sugar Hit

You are here: Home » Federal Politics » Political News

## Research cuts anger universities

October 22, 2012

Bianca Hall
*Political Correspondent*
View more articles from Bianca Hall
Follow Bianca on Twitter

---

# TheScientist
### EXPLORING LIFE, INSPIRING INNOVATION

News ▾   **Magazine** ▾   **Multimedia** ▾   **Subjects** ▾   **Surveys** ▾   **Careers** ▾

Advertisement

Turn gene expression on or turn gene expression off using 365 nm light. **Photo-Morpholinos**

The Scientist » The Nutshell

## Funding Cuts Threaten Big Data
**Reduced support from the US National Library of Medicine threatens to shut down five popular biological databases.**

By Jef Akst | September 5, 2012

---

# nature   *International weekly journal of science*

Search          ▶ Advanced

Home   News & Comment   Research   Careers & Jobs   Current Issue   Archive   Audio & Video   For Author

Archive   Volume 489   Issue 7414   News   Article

*NATURE | NEWS*

## Databases fight funding cuts

Online tools are becoming ever more important to biology, but financial support is unstable.

**Monya Baker**

05 September 2012

---

Fran Berman

# Sustainable Data Infrastructure Starts with a Sustainable Economic Model

**It's not just about the cost of storage**. Research data infrastructure costs increase with usage, stewardship and access requirements, perceived value

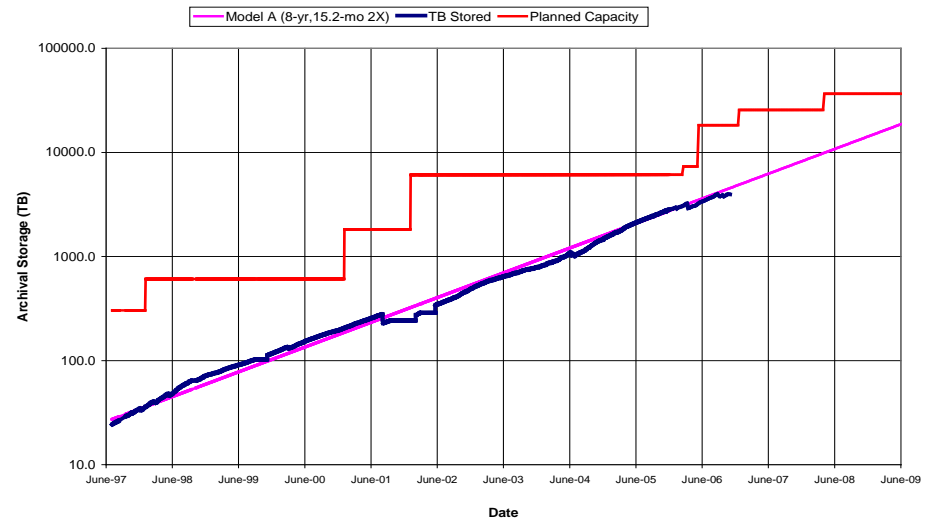**Greater costs beyond the center (including "big" data)**



- More mgt., stewardship required
- Long-lived data
- "Locally Manageable" Data
- Big data
- Coupled data services
- Access control, broad access
- More curation required

Fran Berman

# Data Economics: Data Management, Stewardship, and Use Incur Continuing Infrastructure Costs

## Costs include

- Maintenance and upkeep
- Software tools and packages
- Utilities (power, cooling)
- Space
- Networking
- Security and failover systems
- People (expertise, help, infrastructure management, development)
- Training, documentation
- Monitoring, auditing
- Reporting costs
- Costs of compliance with regulation, etc.

## Resources and Resource Refresh



*SDSC Data Storage Growth '97-'09*

- *Most valuable data replicated*
- *As research collections increase, storage capacity must stay ahead of demand*

*Information courtesy of Richard Moore, SDSC*

Fran Berman

# In the Public Sector, Research and Infrastructure often compete for limited funding.
## Infrastructure Investment a hard sell …

| | Research | Infrastructure |
|---|---|---|
| **What is Newsworthy?** | New discoveries and breakthroughs | Failure of systems |
| **What is the value proposition?** | Domain and national leadership and competitiveness | Enabler of innovation |
| **Funding Model** | Fixed-term funding | Continuous long-term support |



**BBC NEWS SCIENCE & ENVIRONMENT**

8 October 2013 Last updated at 14:01 GMT

Higgs boson scientists win Nobel prize in physics

COMMENTS (643)

By James Morgan
Science reporter, BBC News



Crisp photos of moon landing are missing
Spectacular images of day were stored, forgotten -- and lost

Marc Kaufman, Washington Post
Sunday, February 4, 2007



**nature** International weekly journal of science

Home | News & Comment | Research | Careers & Jobs | Current Issue | Archive | Au

News & Comment > News > 2014 > March > Article

NATURE | NEWS

Scientists losing data at a rapid rate

Decline can mean 80% of data are unavailable after 20 years.

Elizabeth Gibney & Richard Van Noorden

Fran Berman

# Data Infrastructure particularly important in light of increasing National R&D Agency Requirements for Data Access and Management



Fran Berman

# Economics of Public Access:
## Who Pays the Data Bill?



**POLICY**FORUM

SCIENCE PRIORITIES

### Who Will Pay for Public Access to Research Data?

Francine Berman[1] and Vint Cerf[2]

When economic models and infrastructure are not in place to ensure access and preservation, federally funded research data are "at risk."

On 22 February, the U.S. Office of Science and Technology Policy (OSTP) released a memo calling for public access for publications and data resulting from federally sponsored research grants (1). The memo directed federal agencies with more than $100 million R&D expenditures to "develop a plan to support increased public access to the results of research funded by the Federal Government." Perhaps even more succinctly, a subsequent New York Times opinion page sported the headline "We Paid for the Research, So Let's See It" (2). So who pays for data infrastructure?

The OSTP memo requested agencies to provide plans by September 2013 that describe their strategies for providing public access to both research publications and research data. Plans are expected to be implemented using "resources within the existing agency budget," i.e., no new money should be expected. Currently, federal R&D agencies are working hard to foster approaches to public access, to assess needs for supporting partnerships and enabling infrastructure, and to develop timetables and approaches for implementation. We focus here on the research data portion of the OSTP memo.

Research data of community value are supported today in a variety of ways. Some of them, like those in the Protein Data Bank (PDB) (3)—a database of protein structure information used heavily by the life sciences community—are supported by the public sector. (In particular, U.S. funding from the National Science Foundation (NSF), the National Institutes of Health (NIH), and the U.S. Department of Energy for the Research Collaboratory for Structural Bioinformatics (RCSB) PDB is $6.3 million annually.) Other data, as from the Longitudinal Study

What happens to valuable data when project funding ends? Consider, for example, a 3-year research project in which valuable sensor data are collected from an environmentally sensitive area. Those data may be useful not just for the duration of the project but for the next decade or more to collaborators and a broader community of researchers. For the first 3 years, the costs of stewardship (including development of a database that supports analysis, access to the data for the community through a portal, adequate storage and management of the data collection, and so on) may be paid for by the grant. But who pays for subsequent support? In such cases, research data may become more valuable just as the economics of stewardship become less viable.

Up to this point, no one sector has stepped up to take on the problem alone. In the public sector, federal R&D agencies are unrealistic to expect as much. In the to allocate enough resources to support federally funded research data.

*Article: Science Magazine, August 9, 2013. Free public access link at http:/www.cs.rpi.edu/~bermaf/*

# Op-Ed Recommendations:  Partner Across Sectors to Distribute the Preservation and Stewardship Responsibilities

- Facilitate private sector stewardship of public access research data as a public good

- Clarify public sector stewardship commitments: articulate what data will  / won't be supported

**Private Sector**
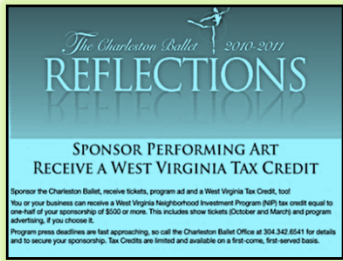
**Public Sector**

**Academia**

**Individuals**

- Create sustainable university library and repository stewardship solutions

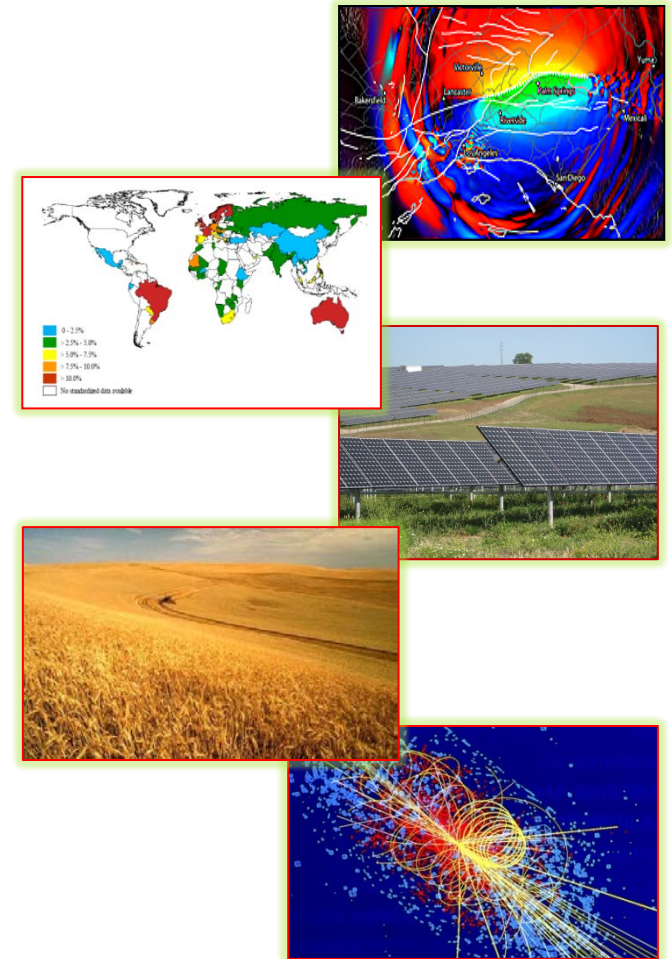- Evolve research culture to take advantage of what works in the private sector



Fran Berman

# Next Steps Towards Sustainable Stewardship

- Identify and evolve an expanding network of repositories for publicly accessible research data

- Create useful metrics for successful stewardship and economic stability that can be used to support the development of effective organizational support

- Create a plan and actionable recommendations for strategic investment

Fran Berman

# Last Words:  Information Infrastructure is necessary for 21st Century Innovation

- **Value Proposition:**

  – Virtually all fields becoming data driven

  – **Adequate and sustainable data infrastructure is critical to drive innovation and HPC applications**

- **What can we do:**

  – Include data stewardship, management, use, access and preservation as part of project planning, budget and efforts

  – Recognize and publish the data contributions of our work as well as the research contributions

**Small steps** (things to do on Monday morning)**:**

1. If you don't have one, create a data management plan for your current project for a reasonable fixed term of time

2. Make your data available to the community (as appropriate) by curating it and ingesting it into a publicly accessible repository

3. Cite and publish your data as appropriate when you write about your results

Fran Berman

# Thank You!



Fran Berman