# Towards a Sustainable Ecosystem for Data Driven Research and Innovation
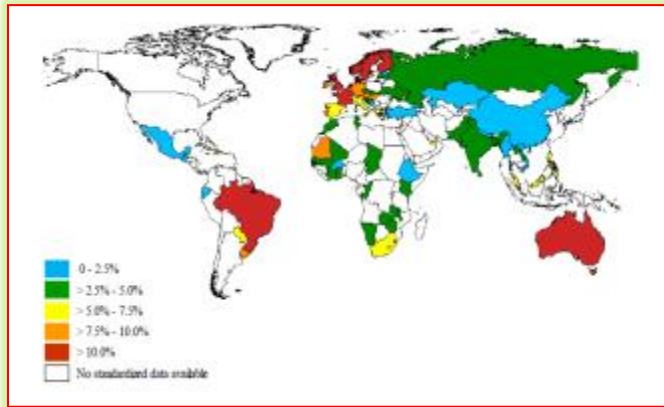
**Dr. Francine Berman**

Chair, Research Data Alliance / US

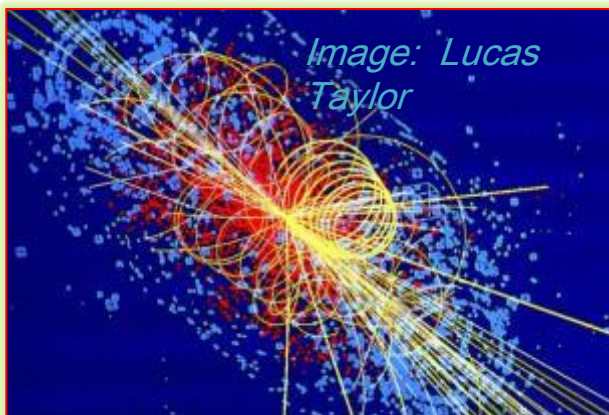Edward P. Hamilton Distinguished Professor of Computer Science, Rensselaer Polytechnic Institute

# Innovation opportunity: Research data driving solutions to 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?



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

Image: Ceinturion, Wikipedia

# Infrastructure reality: Access, use and re-use of data now and in the future presupposes data sustainable stewardship and preservation *today*

- **Stewardship and Preservation critical:** "**Homeless" data cease to exist**

- Sustainable data infrastructure necessary to support

  - Data management plans

  - Public access to research data

  - Use and re-use of data

  - Reproducibility of results





**InformationWeek** CONNECT TECHNOL

Home   News & Commentary   Authors   Slideshows   Video   Reports

STRATEGIC CIO   SOFTWARE   SECURITY   CLOUD   MOBILE

INFRASTRUCTURE // PC & SERVERS

NEWS
11/2/2012
04:19 PM

**Sandy A Grim Reminder: Back Up Your Data**

Once again, disaster — this time Hurricane Sandy — reminds businesses and consumers that they should be backing up their data.
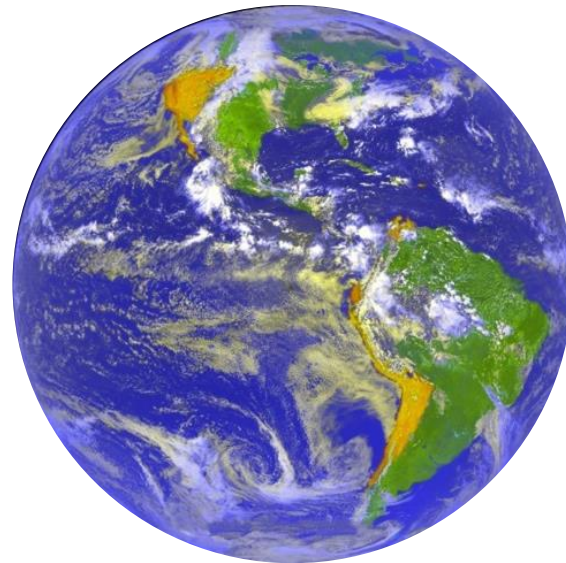
# How can we usefully think about sustainability?

*Sustainable development:* "development that meets the needs of the present without compromising the ability of future generations to meet their own needs."

*Our Common Future, U.N. Brundtland Commission*

- **Key components**
  - Ecological sustainability
  - Cultural / institutional sustainability
  - Economic sustainability
  - Political sustainability

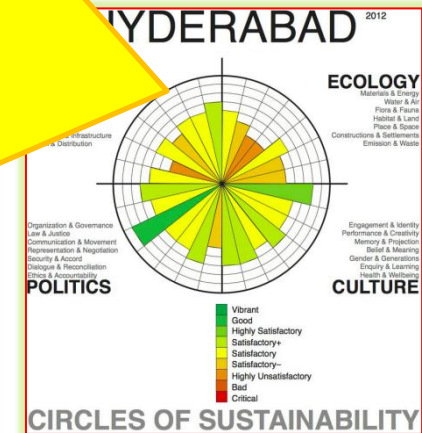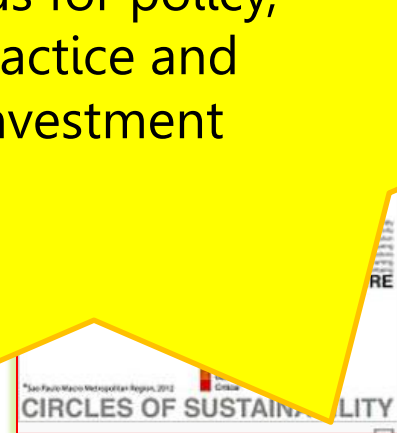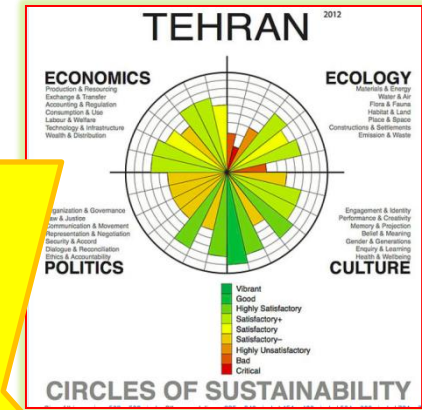# How can we measure sustainability?

- **Circles of Sustainability** developed to assess and understand sustainability. Used

  - For managing proje[...] towards socially sus[...] outcomes

  - To assess the susta[...] cities and urban[...]

- *Used by global orga[...] including the United Nation[...] Compact Cities Programme, T[...] World Association of Metropo[...] World Vision, and others.*
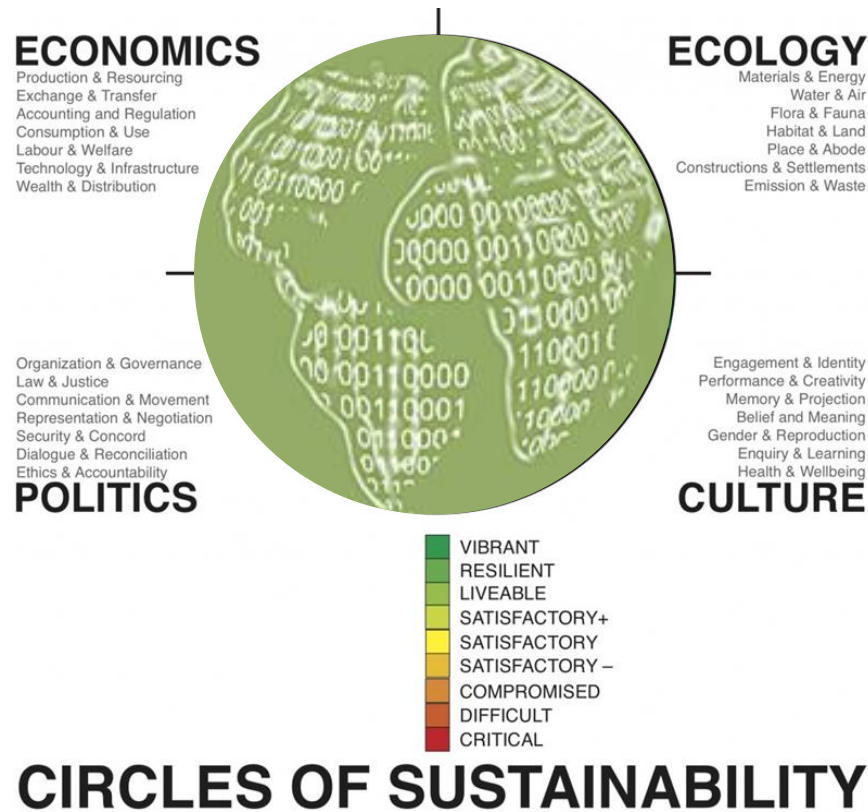
**Assessment → Action**

Sustainability analysis helps identify strategic areas for policy, practice and investment



JOHANNESBURG



TEHRAN 2012

CIRCLES OF SUSTAINABILITY



HYDERABAD 2012

CIRCLES OF SUSTAINABILITY

CIRCLES OF SUSTAINABILITY

# Circles of sustainability provide a useful lens with which to consider sustainability of digital research data

*How can we create a viable support model for digital access and preservation?*

*Why is digital preservation and access such a hard sell?*



**ECONOMICS**
Production & Resourcing
Exchange & Transfer
Accounting and Regulation
Consumption & Use
Labour & Welfare
Technology & Infrastructure
Wealth & Distribution

Organization & Governance
Law & Justice
Communication & Movement
Representation & Negotiation
Security & Concord
Dialogue & Reconciliation
Ethics & Accountability
**POLITICS**

**ECOLOGY**
Materials & Energy
Water & Air
Flora & Fauna
Habitat & Land
Place & Abode
Constructions & Settlements
Emission & Waste

Engagement & Identity
Performance & Creativity
Memory & Projection
Belief and Meaning
Gender & Reproduction
Enquiry & Learning
Health & Wellbeing
**CULTURE**

VIBRANT
RESILIENT
LIVEABLE
SATISFACTORY+
SATISFACTORY
SATISFACTORY −
COMPROMISED
DIFFICULT
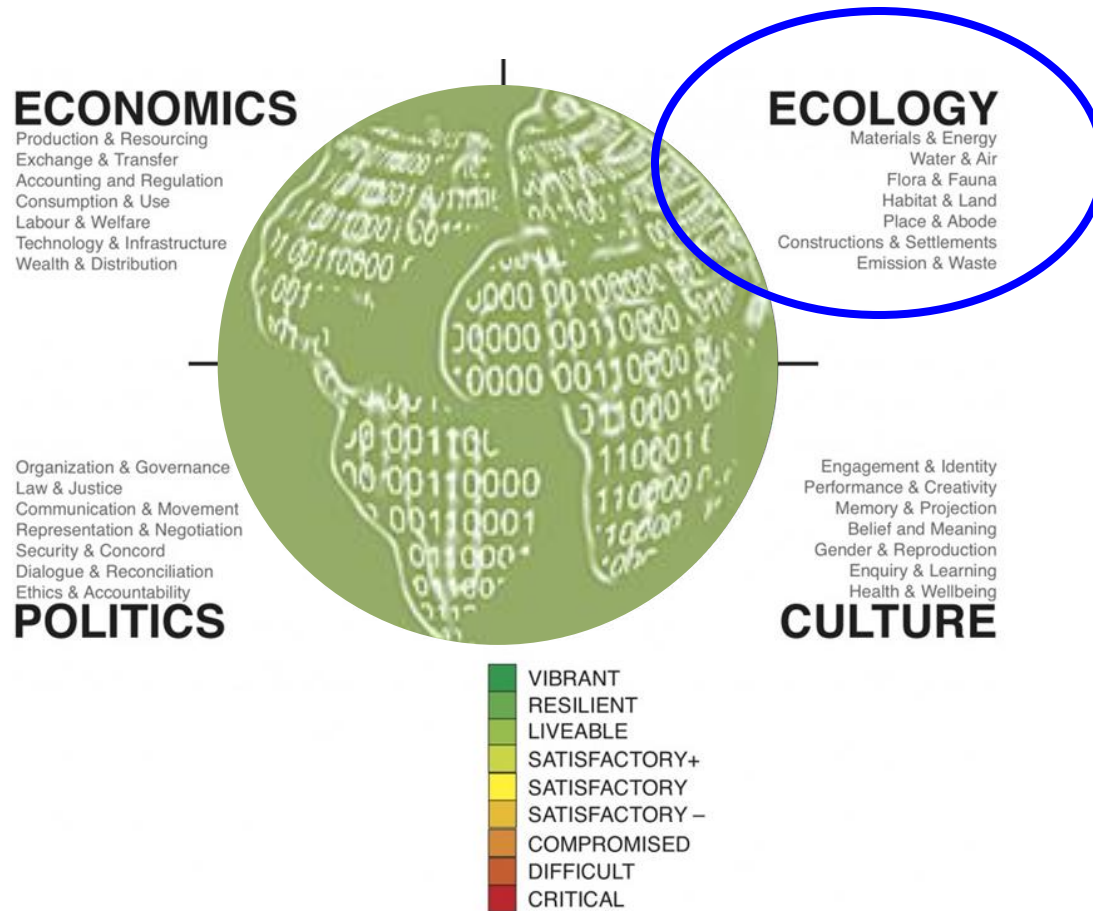CRITICAL

## CIRCLES OF SUSTAINABILITY

*What infrastructure is needed to support digital stewardship and preservation?*

*How do we maximize the access, sharing and exchange of digital research data?*

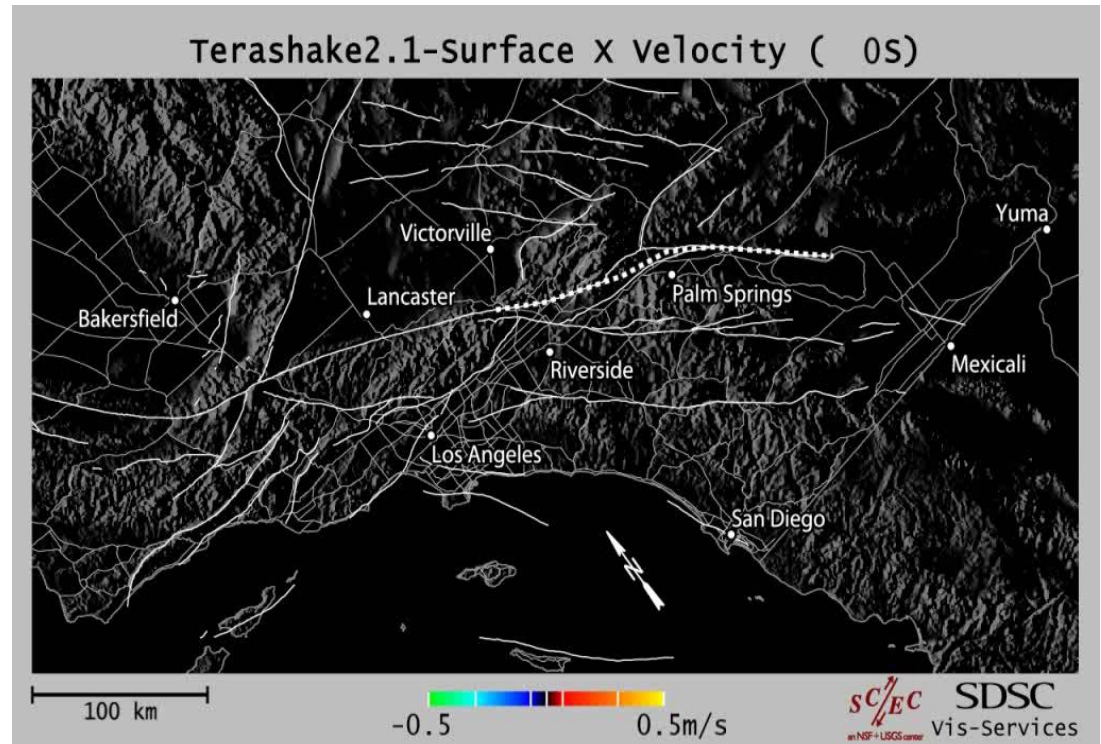# What infrastructure is needed to support digital stewardship and preservation?



CIRCLES OF SUSTAINABILITY

# 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*

# 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*

# Terashake data infrastructure*

- **Data Management**
  - 10 Terabytes moved per day during execution over 5 days
  - Derived data products registered into SCEC digital library (total SCEC library has 168 TB)

- **Data Post-processing:**
  - *Movies* of seismic wave propagation
  - Seismogram formatting for interactive on-line analysis
  - *Derived data:*
    - Velocity magnitude
    - Displacement vector field
    - Cumulative peak maps
    - Statistics used in visualizations

\* circa ~2007

# TeraShake Resources

## Computers and Systems

- 80,000 hours on IBM Power 4 (DataStar)
- 256 GB memory p690 used for testing, p655s used for production run, TeraGrid used for porting
- 30 TB Global Parallel file GPFS
- Run-time 100 MB/s data transfer from GPFS to SAM-QFS
- 27,000 hours post-processing for high resolution rendering

## People

- 20+ people for IT support
- 20+ people in domain research

## Storage

- SAM-QFS archival storage
- HPSS backup
- Storage Resource Broker collection with 1,000,000 files

# Technical infrastructure critical part of the ecosystem for data-driven innovation

Data access via **portals**, science gateways, etc.

Database and data collection **systems**

Data **services** to support use and re-use

Data analysis **algorithms**, data-driven models and simulations

Data **visualization** tools

Semantic **frameworks**

Data **management** systems

Data **storage**

# Social, organizational, and human infrastructure for data-driven results equally important

Policy

Systems Interoperability

Common Standards

Sustainable Economics

Community Practice

Workforce and training

**Fran Berman**          12

# Data-driven ecosystem requires multiple kinds of infrastructure

Data Stewardship Economic Support

Data-driven Innovation

## Data Infrastructure

**Technical Infrastructure**
SW and systems, Tools and algorithms, Hardware and facilities

**Social Infrastructure**
Policy, Practice, Standards, Rights, Community culture

# How can we create a viable support model for digital access and preservation?



**ECONOMICS**
Production & Resourcing
Exchange & Transfer
Accounting and Regulation
Consumption & Use
Labour & Welfare
Technology & Infrastructure
Wealth & Distribution

**ECOLOGY**
Materials & Energy
Water & Air
Flora & Fauna
Habitat & Land
Place & Abode
Constructions & Settlements
Emission & Waste

Organization & Governance
Law & Justice
Communication & Movement
Representation & Negotiation
Security & Concord
Dialogue & Reconciliation
Ethics & Accountability
**POLITICS**

Engagement & Identity
Performance & Creativity
Memory & Projection
Belief and Meaning
Gender & Reproduction
Enquiry & Learning
Health & Wellbeing
**CULTURE**

VIBRANT
RESILIENT
LIVEABLE
SATISFACTORY+
SATISFACTORY
SATISFACTORY –
COMPROMISED
DIFFICULT
CRITICAL

**CIRCLES OF SUSTAINABILITY**

# Data economics: Responsible data stewardship requires a viable business model for sustaining its underlying infrastructure

**Data infrastructure costs increase** with usage, stewardship and access requirements, perceived value

**Greater costs at the extremes (including "big" data) …**

# It's not just about the cost of storage

## Data Infrastructure 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, policy, 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*

# Increased requirements for access mean increased need for data infrastructure

- **Increasing U.S. ... requirements ...** research data

- **February 2013 ...** for public access ... publications

  - Strategy for ca... public-private...

  - Strategy for in... access, dissem...

  - **No new mone...** the existing a...

# Economics of public access:
## Who pays the data bill?



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

# Op-ed recommendations:  Cultivate / coordinate preservation and stewardship options in every sector



Private Sector

Public Sector

Academia

Individuals

# Op-ed recommendations: Cultivate / coordinate preservation and stewardship options in every sector

## Private Sector

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

## Public Sector

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

Not govt. supported **??** Govt. supported

# Op-ed recommendations:  Cultivate / coordinate preservation and stewardship options in every sector

## Academic Sector

- Create sustainable university library and repository stewardship solutions



## Individuals

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

# No magic economic bullet. Coordination between approaches can provide even more robust options for stewardship

# Why is digital preservation and access such a hard sell?

# Academic and public sector infrastructure challenges

| | 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 |
| **What is the funding model?** | Fixed-term funding | Continuous long-term support |
| **Who is responsible?** | Various govt. R&D agencies, NGOs, etc. | No-one's major priority |



Stephanie A. Miner, the Syracuse mayor, said [infrastructure is] too often overlooked when politicians want to spend money on economic development. **"You don't cut ribbons for new water mains, but that's really what matters."**
NY Times, Feburary 15, 2014

# Systemic challenges to sustainable stewardship

*(from the Blue Ribbon Task Force Interim Report [at brtf.sdsc.edu])*

- **Poor alignment between stakeholders** in the digital preservation and access world and their roles, responsibilities and support models

- There is a **lack of institutional, enterprise, and/or community incentives** to support the collaboration needed to enforce sustainable economic models

- **Complacency that current practices are "good enough"** and / or the **problem is not urgent.** Both "carrots" (in the form of recognition that access to information is an investment in current and future success) and "sticks" (in the form of penalties for non-compliance, accounting of explicit opportunity costs, or costs of lost information) are needed

- Fear that digital access and **preservation is too big to take on**

# Reports are not enough, but good reports can provide compelling evidence needed by stakeholders for action

# Making the case: Quantifying / qualifying advancement and fear

- **Political capital critical for prioritization and investment in digital stewardship and preservation**

- *Arguments that influence stakeholders and their enablers:*

  – Better economic growth / more jobs

  – Greater leadership / accelerated innovation

  – Increased reputation / competitive advantage

  – Fear of disaster / loss of reputation



Big data—capturing its value

**$300 billion**
potential annual value to US health care—more than double the total annual health care spending in Spain

**€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

**Big Data: The next frontier for innovation, competition and productivity**
*McKinsey Global Institute, 2011*



Forbes — Target Shares Tumble As Retailer Reveals Cost Of Data Breach

RESEARCH DATA ALLIANCE

# How do we maximize the access, sharing and exchange of digital research data?



ECONOMICS
Production & Resourcing
Exchange & Transfer
Accounting and Regulation
Consumption & Use
Labour & Welfare
Technology & Infrastructure
Wealth & Distribution

ECOLOGY
Materials & Energy
Water & Air
Flora & Fauna
Habitat & Land
Place & Abode
Constructions & Settlements
Emission & Waste

Organization & Governance
Law & Justice
Communication & Movement
Representation & Negotiation
Security & Concord
Dialogue & Reconciliation
Ethics & Accountability
**POLITICS**

Engagement & Identity
Performance & Creativity
Memory & Projection
Belief and Meaning
Gender & Reproduction
Enquiry & Learning
Health & Wellbeing
**CULTURE**

VIBRANT
RESILIENT
LIVEABLE
SATISFACTORY+
SATISFACTORY
SATISFACTORY –
COMPROMISED
DIFFICULT
CRITICAL

**CIRCLES OF SUSTAINABILITY**

# Data-sharing driving innovation across research cultures

# World-wide national and professional communities focusing on 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

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**



**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 |

RESEARCH DATA ALLIANCE

# The Research Data Alliance (RDA)

- Global community-driven organization launched in March 2013 to accelerate data-driven innovation

- RDA focus is on building the **social, organizational and technical infrastructure** to

  - *reduce barriers to **data sharing and exchange***

  - *accelerate the development of coordinated global data infrastructure*

# 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**

# The RDA Community Today: Over 2300 members from 96 countries (as of 9/14)



Africa 3%
Asia 5%
South America 1%
Austral-pacific 5%
North Amer. 36%
Eur. 50%

Other 6%
Private 13%
Govt. 18%
Acad. 63%

Map courtesy traveltip.org

Fran Berman          33

# Precipitous growth


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

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

First Organizational Partner Meet-up

First BOFs

First Organizational Assembly

6 co-located events

14 BOF, 12 Working Groups, 22 Interest Groups


**RDA Plenary 3**
**Dublin, Ireland**

First Working Groups and Interest Groups

**240 participants**

380 participants from 22 countries

497 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


**RDA Plenary 2**
**Washington, DC**

First Working Group exchange meeting


**RDA Plenary 4**
Amsterdam

# RDA Interest (IG) and Working Groups (WG) by focus 1 (as of 9/14)

*under review

## Domain Science - focused

- Toxicogenomics Interoperability IG
- Structural Biology IG
- Biodiversity Data Integration IG
- Agricultural Data Interoperability IG
- Wheat Data Interoperability WG
- Digital Practices in History and Ethnography IG
- Geospatial IG

- Marine Data Harmonization IG
- Metabolomics IG
- RDA/CODATA Materials Data Infrastructure and Interoperability IG
- Research Data Needs of the Photon and Neutron Science Community IG
- Defining Urban Data Exchange for Science IG*
- The BioSharing Registry:  Connecting data policies, standards and databases in the life sciences WG*
- Urban Quality of Life Indicators WG*

## Community Needs - focused

- Community Capability Model IG
- Engagement IG
- RDA / CODATA Summer Schools in Data Science and Cloud Computing in the Developing World WG*

- Development of Cloud Computing Capacity and Education in Developing World Research IG
- Data for Development IG
- Education and Training on handling of research data IG

# RDA Interest (IG) and Working Groups (WG) by focus 2 (as of 9/14)

## Reference and Sharing - focused

- Data Citation WG
- Standardization of Data Categories and Codes WG
- RDA/CODATA Legal Interoperability IG
- Reproducibility IG*
- Data Description Registry Interoperability Working Group
- RDA / WDS Publishing Data Bibliometrics WG

## Data Stewardship and Services - focused

- Research Data Provenance IG
- Preservation e-infrastructure IG
- RDA / WDS Publishing Data Services WG
- RDA / WDS Publishing Data Workflows WG
- Long-tail of Research Data IG
- RDA/WDS Publishing Data IG
- RDA/WDS Repository Audit and Certification WG
- Domain Repositories Interest Group
- Brokering Interest Group
- ELIXIR Bridging Force IG*
- Libraries for Research Data IG*
- RDA / WDS Certification of Digital Repositories IG
- RDA / WDS Publishing Data Cost Recovery for Data Centres IG

## Base Infrastructure - focused

- Data Foundation and Terminology WG
- Metadata Standards Directory WG
- Practical Policy WG
- PID Information Types WG
- Data Type Registries WG
- Data in Context IG
- Big Data Analytics IG
- Data Brokering WG*
- Federated Identity Management IG
- Metadata IG
- PID Interest Group
- Service Management IG
- Data Fabric IG

# Organizations committed to joining RDA

- **Organizational Members:**
  - Alliance for Permanent Access
  - American University Library
  - Australian National Data Service
  - Barcelona Supercomputing Center - Centro Nacional de Supercomputación
  - Columbia University Library
  - CNRI
  - CSC
  - Digital Curation Center
  - EIROForum IT Working Group
  - eResearch Services and Scholarly Application Development Division of Information Services, Griffith University
  - European Data Infrastructure (EUDAT)
  - National Institute of Advanced Industrial Science and Technology (AIST), Japan
  - International Association of STM Publishers
  - Internet2

  - Microsoft Research
  - NZ eScience Infrastructure
  - Purdue University Libraries
  - Research Data Canada
  - Scholarly Publishing and Academic Resources Coalition (SPARC)
  - Washington University in St. Louis Libraries
  - Science and Technology Facilities Council

- **Affiliates**
  - CODATA
  - ICSU World Data System
  - ORCID
  - DataCite
  - Global Alliance for Genomics and Health
  - CASRAI

# Accelerating data sharing infrastructure, coalescing culture: 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
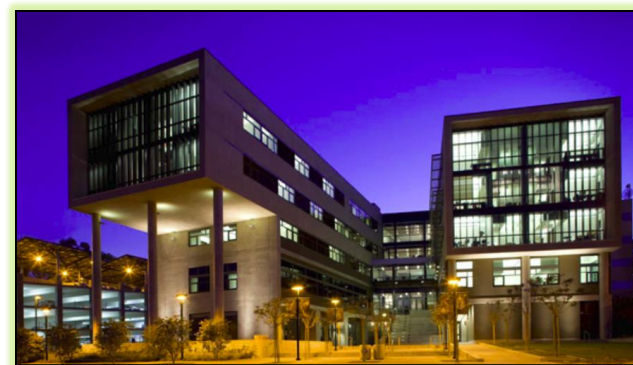
**Partnership with Organizations**

New partners in the Organizational Assembly

Focused strategy to support development of industry infrastructure for data sharing

# 2015 Focus: Development and adoption of first RDA deliverables



- **March 9-11, 2015: RDA/US hosting RDA Plenary 5 in San Diego**

    – Working Meeting of the RDA with co-located community meetings

- **San Diego Supercomputer Center hosting RDA Adoption Day on March 8**

# Thank You

## ECONOMICS

**Budget realistically** for the costs of data stewardship and preservation

**Prioritize the "data bill"** at the same level as other critical infrastructure.

## ECOLOGY

Create and implement a **data management and stewardship plan** for your project for a reasonable fixed term of time.

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

## POLITICS

**Adopt / support policy and practice** that enables the development and continued maintenance of sustainable stewardship, data sharing, and broad access

## CULTURE

Contribute /create a local  / community culture of **data sharing**

**Cite and publish your data** when you write about your results.  Work with your professional societies and conferences to include **"data sessions"** *(idea from Sibel Adali)*