

Leadership Team: William Michener¹, Suzie Allard², John Cobb³, Robert Cook³, Patricia Cruse⁴, Mike Frame⁵, Stephanie Hampton⁶, Vivian Hutchison⁷, Matthew Jones⁸, Steve Kelling⁷, Rebecca Koskela¹, Carol Tenopir², Dave Vieglais⁸, Todd Vision⁹, Bruce Wilson²

Co-Investigators: Paul Allen⁷, Peter Buneman¹⁰, Randy Butler¹¹, Ewa Deelman¹², David DeRoure¹³, Cliff Duke¹⁴, Carole Goble¹⁵, Donald Hobern¹⁶, Peter Honeyman¹⁷, Jeffery Horsburgh¹⁸, John Kunze⁴, Bertram Ludascher¹⁹, Maribeth Manoff², Line Pouchard², Robert Sandusky²⁰, Ryan Scherler⁹, Mark Servilla¹, Jake Weltzin¹

¹University of New Mexico; ²University of Tennessee; ³Oak Ridge National Laboratory; ⁴University of California - California Digital Library; ⁵U.S. Geological Survey; ⁶National Center for Ecological Analysis and Synthesis - University of California - Santa Barbara; ⁷Cornell University; ⁸University of Kansas; ⁹National Evolutionary Synthesis Center, University of North Carolina; ¹⁰University of Edinburgh; ¹¹University of Illinois - Urbana Champaign; ¹²University of Southern California; ¹³University of Southampton; ¹⁴Ecological Society of America; ¹⁵University of Manchester; ¹⁶Atlas of Living Australia; ¹⁷University of Michigan; ¹⁸Utah State University; ¹⁹University of California - Davis; ²⁰University of Illinois - Chicago

Abstract:
Addressing the Earth's environmental problems requires that we change the ways that we do science; harness the enormity of existing data; develop new methods to combine, analyze, and visualize diverse data resources; create new, long-lasting cyberinfrastructure; and re-envision many of our longstanding institutions. DataONE (Observation Network for Earth) represents a new virtual organization whose goal is to enable new science and knowledge creation through universal access to data about life on earth and the environment that sustains it.

DataONE is designed to be the foundation for new innovative environmental science through a distributed framework and sustainable cyberinfrastructure that meets the needs of science and society for open, persistent, robust, and secure access to well-described and easily discovered Earth observational data.

Supported by the U.S. National Science Foundation, DataONE will ensure the preservation and access to multi-scale, multi-discipline, and multi-national science data. DataONE is transdisciplinary, making biological

data available from the genome to the ecosystem; making environmental data available from atmospheric, ecological, hydrological, and oceanographic sources; providing secure and long-term preservation and access; and engaging scientists, land-managers, policy makers, students, educators, and the public through logical access and intuitive visualizations. Most importantly, DataONE will serve a broader range of science domains both directly and through the interoperability with the DataONE distributed network. DataONE is a five year project that began in Fall 2009 (William Michener, PI, University of New Mexico).

The Vision: "DataONE will be commonly used by researchers, educators, and the public to better understand and conserve life on earth and the environment that sustains it."

By creating an infrastructure of technology and standards, people, and institutions to support the full life cycle of biological, ecological, and environmental data and tools that enable universal access, DataONE will accelerate use of earth observational data in research, education and decision-making. In so doing, DataONE will transform our understanding of ecological processes and conserve life on earth and the environment that sustains it.

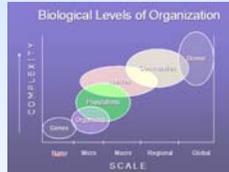
Why do we need DataONE?: Societal and Environmental challenges



Popular press and results from the International Geosphere Biosphere Program show that environmental challenges are of increasing concern for us all.

Science Challenges

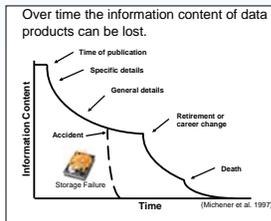
- To understand, we need easy access to different types of data
- Wide range of spatial and temporal scales (plot data to remote sensing data)
- Breadth of science domains (biological, environmental, social, and economic)
- Citizen science networks, of increasing importance



Data Challenges

Scientists need access to the data generated by research to verify findings and test new hypotheses.

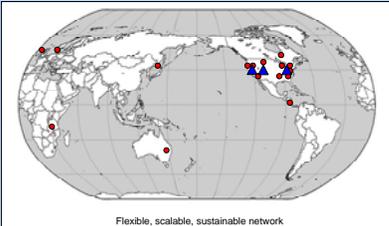
- Poor data practices place the scientific record at risk
 - Data are massively dispersed
 - Not on the Web
 - Orphaned
 - Multiple Semantics
 - Not easily discovered
 - Poor data practices
 - Documentation
 - Formats / obsolescence
 - Lack of Standards
 - Poor stewardship
 - Media Obsolescence
 - Heterogeneous, incompatible formats
 - Difficult to combine data from diverse sources



What is DataONE?

Cyberinfrastructure enterprise (tools, ideas, people)

Distributed Framework of Coordinating Nodes and Member Nodes



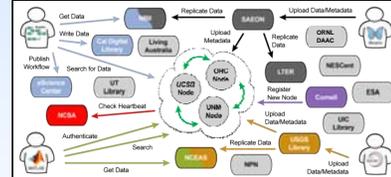
- Coordinating Nodes**
 - retain complete metadata catalog
 - perform basic indexing
 - provide network-wide services
 - ensure data availability (preservation)
 - provide replication services
- 1. University of New Mexico
- 2. University of California-Santa Barbara
- 3. Oak Ridge Campus (UT & ORNL)

- Member Nodes**
 - diverse institutions
 - serve local community
 - provide resources for managing their data
- Member Nodes Deployed in 2010:
 - ORNL DAAC, Knowledge Network for Biocomplexity, Dryad, USGS NBII, and California Digital Library

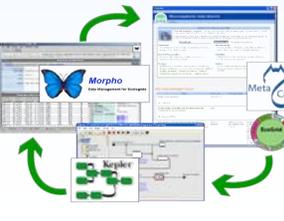
Provide tools to better manage data

Provide support for the entire data life cycle—preparation of data sets, stewardship, tools to access and use the data

- Collection/Preparation
- Deposition/acquisition/ingest
- Curation and metadata management
- Protection, including privacy
- Discovery, access, use, and dissemination
- Interoperability, standards, and integration
- Exploration, visualization, and analysis



Metadata creation, management
Search catalog
Work Flows



Metadata catalog with >70,000 records

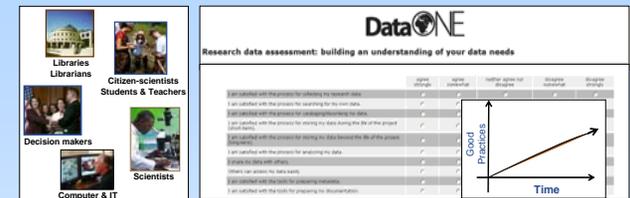


How will we build DataONE?

Understand the community's need, how the community envision solutions

Perform baseline and iterative Community Assessments and Usability Studies

To see where data practices and policies are now, so we can see how practices change over the life of DataONE:



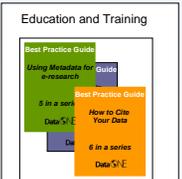
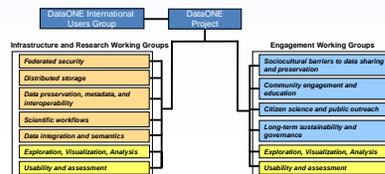
DataONE Assessment: <http://ovvici.com/wsb.dll/e/aaeg3c1e6>

Leverage existing CI whenever possible, build new CI whenever necessary

- Many existing open source efforts exist
 - Metadata Editors: *Mercury, Morpho*
 - Data management: *MATT, UDig, Specify*
 - Analysis and modeling: *R, Octave, netCDF*
 - Workflow systems: *Kepler, Taverna, VisTrails*
 - Grid systems: *Condor, Globus, BOINC*
 - Data and workflow portals: *VegBank, myExperiment*
- Commercial tools important too
 - MATLAB, SAS, ArcGIS, R*
- DataONE: help communities build their own tools
 - Integrate, interoperate, stabilize
 - Create libraries to DataONE Service Interface

Engage the community, reach out, educate, enable new science, and demonstrate success

- DataONE will use Working Groups
- | Structure | Success |
|--|---|
| <ul style="list-style-type: none"> 10 - 20 participants Deep analysis Intensive collaboration Neutral territory 2 week-long meetings per year | <ul style="list-style-type: none"> Community-driven Inclusive High productivity High impact |



- Career Long Learning:
- best practice guides
 - exemplary data management plans
 - podcasts, web-casts
 - workshops and seminars
 - downloadable curricula