

Toward Integrated Environmental Modeling Using Research Data Infrastructures

Jeffery S. Horsburgh

Civil and Environmental Engineering
Utah State University

The Motivators

(from the perspective of an Environmental Engineer interested in water)

- Uncertain climate, land use change, population growth, will alter hydrology, water availability, and water quality
- We need better information to manage water resources
 - But the questions are now transdisciplinary
 - The next generation of models will look very different
- Significant challenges in integrated modeling
 - Coupling models/model components from different disciplines
 - Matching data to required model inputs and outputs at spatial and temporal scales relevant to decision making

Research Data Infrastructures

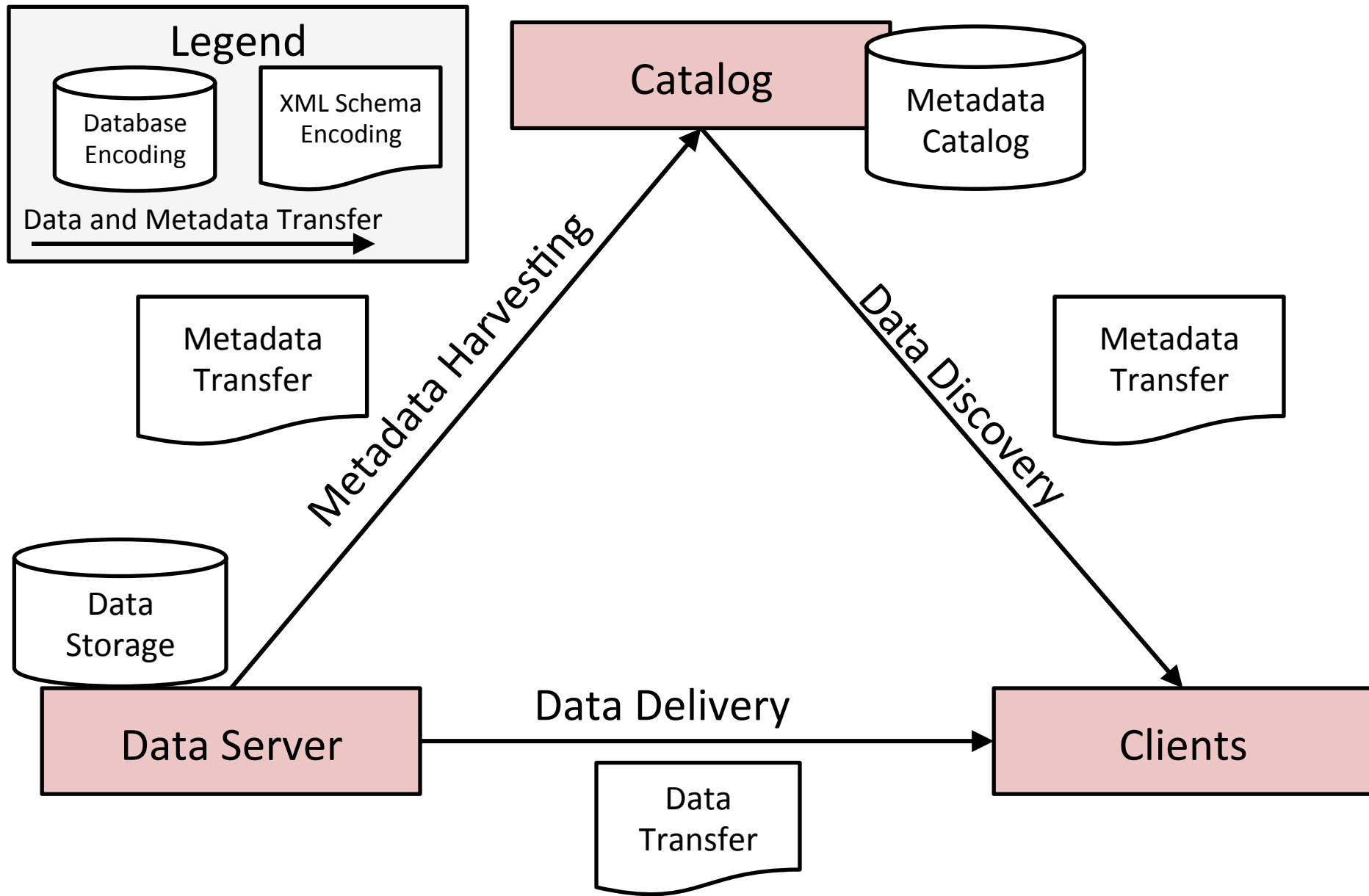
- Research data infrastructures in the U.S. are still disciplinary:
 - Hydrology Community: CUAHSI Hydrologic Information System and Water Data Center
 - Ecological Community: DataONE, Knowledge Network for Biocomplexity (KNB)
 - Critical Zone Community: Critical Zone Observatory Integrated Data System (CZOData)
 - Ocean Community: Integrated Ocean Observing System (IOOS)
 - ...

And, it's not uncommon to find hydrologic data in all of these systems!

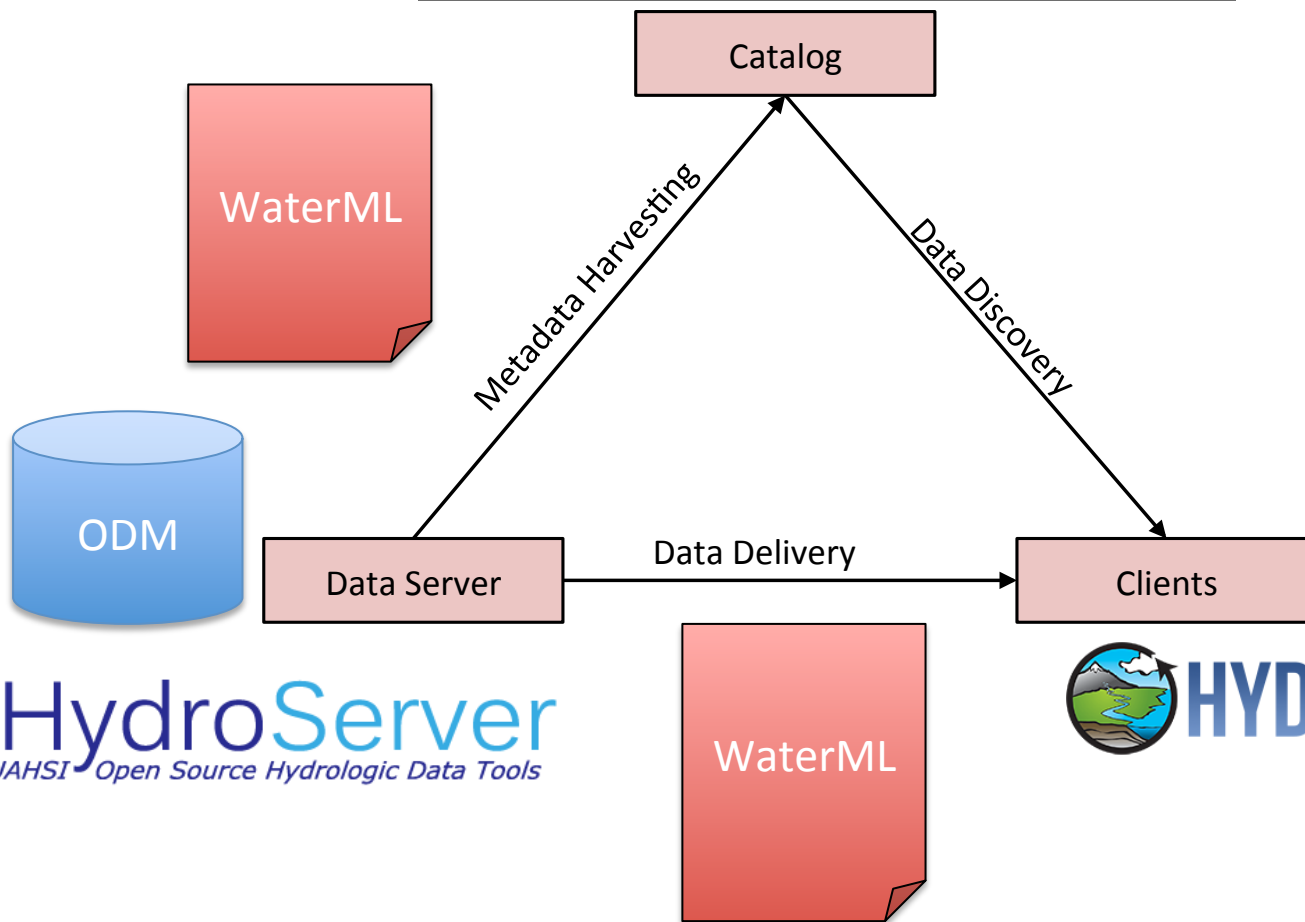
Research Data Infrastructures

- Built using principles of service-oriented architecture (SOA)
- Rely on standard data encodings – but domain specific
- In some cases rely on standard semantics – but not across systems and domains
- Focused on publishing or sharing data on the Internet
 - Common functions: metadata harvest, catalog, search/discovery, data access

Services Oriented Architecture



CUAHSI Water Data Center

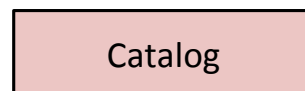
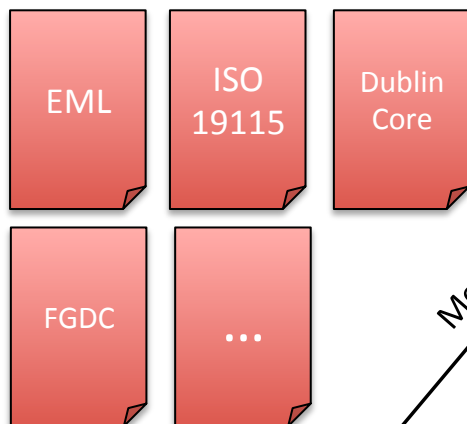


HYDRODESKTOP



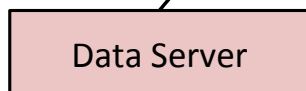
Coordinating Nodes

Multiple metadata standards

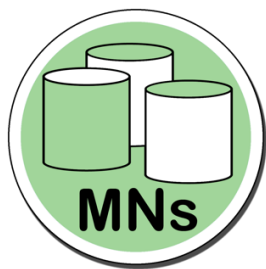
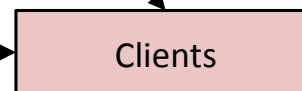


Metadata Harvesting

Data Discovery

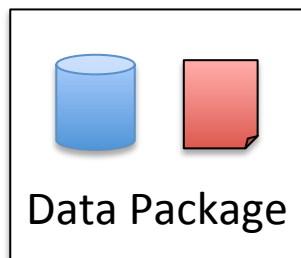


Data Delivery



Member Nodes

No restrictions on data format(s)



Investigator Toolkit

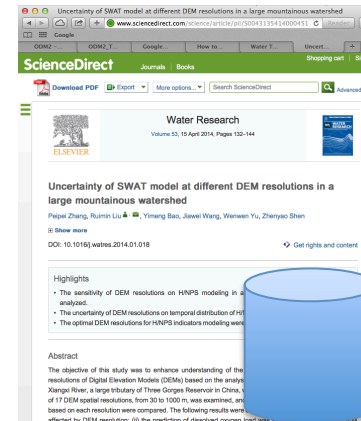


Reproducible Science Versus Integrated Modeling

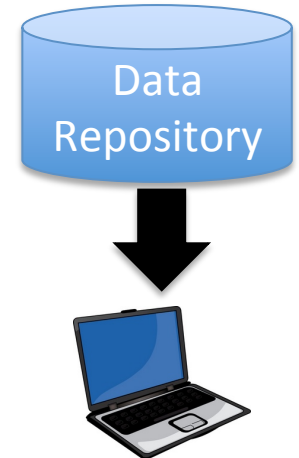
Non-reproducible
science
using integrated
modeling
using SWAT



Search for model
application papers



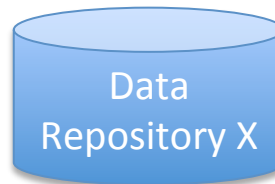
Paper cites data
source



Download data and
run the model

Reproducible Science Versus Integrated Modeling

Now I think
reproducible science
is the only way
to get the junk
out of the data.



Query the data
repository for
“Streamflow”



Query returns
1000 metadata
results



Reproducible Science Versus Integrated Modeling

- Sharing and publication is important for reproducible science
 - Supported well by existing research data infrastructures
- But, the next scientist/modeler looking for data to perform a different analysis
 - Has trouble finding what they need
 - Data are site/study specific and “outsiders” struggle to determine if the data are appropriate for their use
 - Many models require data at spatial and temporal scales not found in research data infrastructures

Challenges

- It takes a “data savvy” investigator considerable effort to discover and access datasets from multiple repositories for a synthetic analysis
- There is a shortage of “data savvy” investigators
- There is currently a mismatch between the functionality of research data infrastructures and what is required for integrated modeling

A Hypothetical Query

I am interested in modeling the effects of aquatic nutrient concentrations on stream metabolism.

“Show me locations where high frequency stream discharge, water temperature, and dissolved oxygen data have been collected in third order streams for which samples of nitrogen and phosphorus have also been collected and that are within one mile of a weather station that measures solar radiation for the same time period.”

Spatial Context

“Show me **locations** where high frequency stream discharge, water temperature, and dissolved oxygen data have been collected **in third order streams** for which samples of nitrogen and phosphorus have been collected and that are **within one mile of a weather station** that measures solar radiation for the same time period.”

Temporal Context

“Show me locations where **high frequency** stream discharge, water temperature, and dissolved oxygen data have been collected in third order streams for which samples of nitrogen and phosphorus have been collected and that are within one mile of a weather station that measures solar radiation **for the same time period.**”

Measured Variables

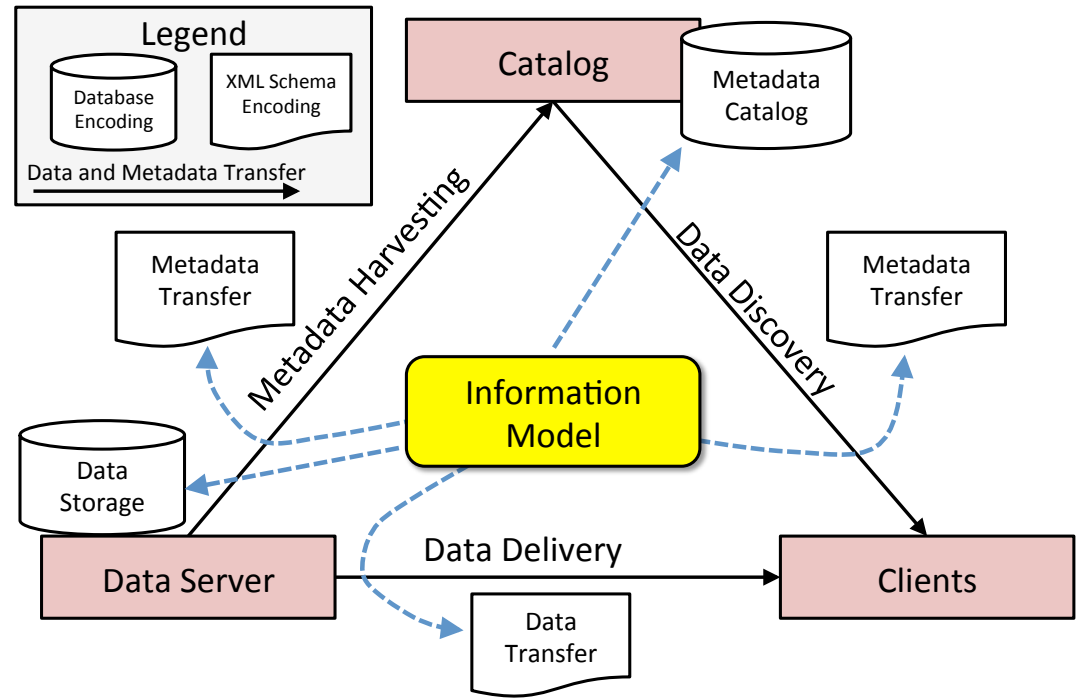
“Show me locations where high frequency **stream discharge, water temperature, and dissolved oxygen data** have been collected in third order streams for which samples of **nitrogen and phosphorus** have been collected and that are within one mile of a weather station that measures **solar radiation** for the same time period.”

Toward Integrated Environmental Modeling Using Research Data Infrastructures

- A vision for next generation research data infrastructures
 - An integrated view of the best available “data” that holistically represent the earth’s hydrologic and environmental systems
 - Seamless presentation of and access to data across existing earth science data repositories
- What is needed:
 - Common information modeling
 - Linking data to the geo-environment
 - Capturing the knowledge content of data
 - Integration with government and agency repositories

Common Information Modeling

- Challenges:
 - Geoscience cyberinfrastructures represent common informational elements inconsistently
 - Opportunities for standardized representation of **spatial, temporal, and measured variable** contexts
 - Semantic and syntactic heterogeneity are major hurdles!



Information Model: Agreement about information needed to describe observational data

Observation Information Models

OGC's Observations & Measurements

OGC and ISO 19156:2011(E)

Open Geospatial Consortium

Approval Date: 2011-12-12

Publication Date: 2013-09-17

External OGC identifier: <http://www.opengis.net/doc/is/om/2.0>

Reference number of this document: OGC 10-004r3

Version: 2.0

Category: OGC® Standard: Abstract Specification

Editors: Simon Cox

OGC Abstract Specification

Geographic information — Observations and measurements

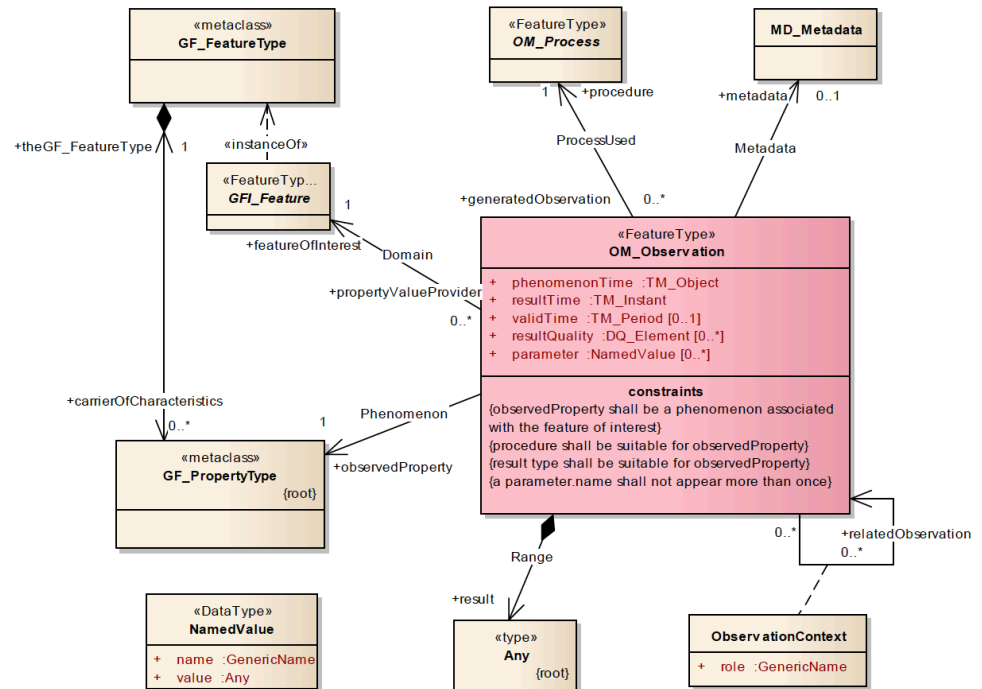
Copyright © 2013 Open Geospatial Consortium
To obtain additional rights of use, visit <http://www.opengeospatial.org/legal/>.

Warning

This document is an OGC Member approved international standard. This document is available on a royalty free, non-discriminatory basis. Recipients of this document are invited to submit, with their comments, notification of any relevant patent rights of which they are aware and to provide supporting documentation. This standard was jointly developed between the OGC and ISO TC 211 and is double branded.

Document type: OGC® Abstract Specification
Document subtype: Encoding
Document stage: Approved for Public Release
Document language: English

Copyright © 2013 Open Geospatial Consortium



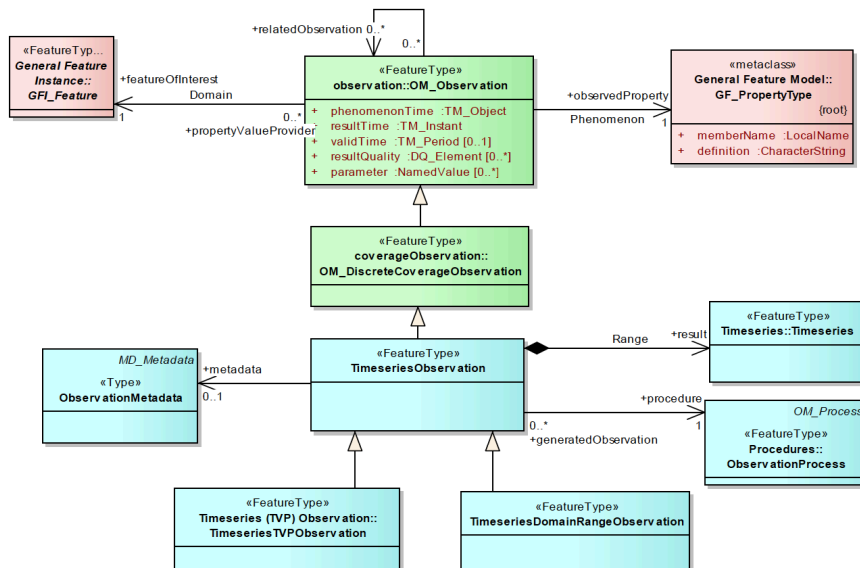
Observation

- Sampling Feature
- Feature of Interest
- Procedure
- Phenomenon/observed property
- Result

In Practice: Profile!

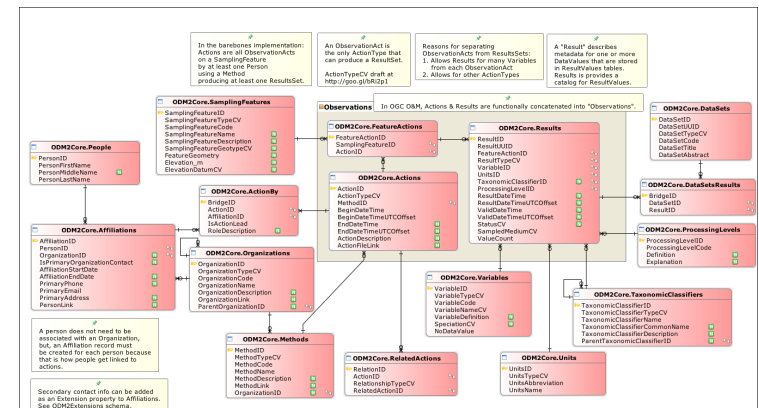
WaterML 2.0

- International standard for time series of water observations data

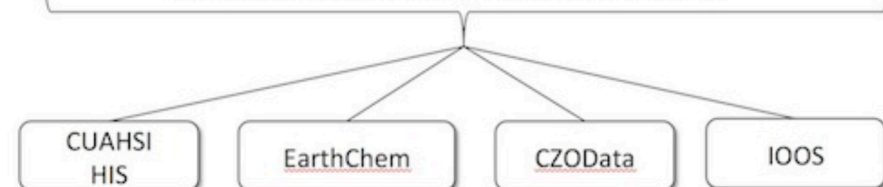


ODM2

- Additional result types beyond time series
- Water quality and solid earth samples, sections, transects, profiles



Common Semantics for Earth Observations



Linking Data to the Geo-Environment

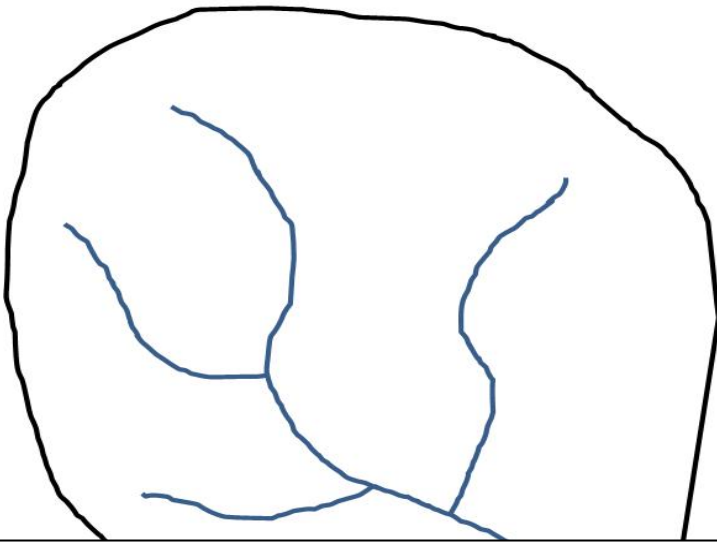
(aka – Spatial Context)

- Feature of Interest – a real world feature whose properties are observed
- Geospatial support – the area or volume to which an observation applies
- Geospatial context – the relationships between a feature of interest and its surrounding features

Linking Data to the Geo-Environment

An Example from Hydrology

Geographic Features



Currently, data users must determine that the gage lies on a particular river, measures the outflow of a particular catchment, is downstream of another gage, shares its location with a water quality monitoring site, and is located near a weather station.

Observations

1	1/1/2012	Streamflow	100cfs
2	1/2/2012	Streamflow	101cfs
3	1/3/2012	Streamflow	102cfs

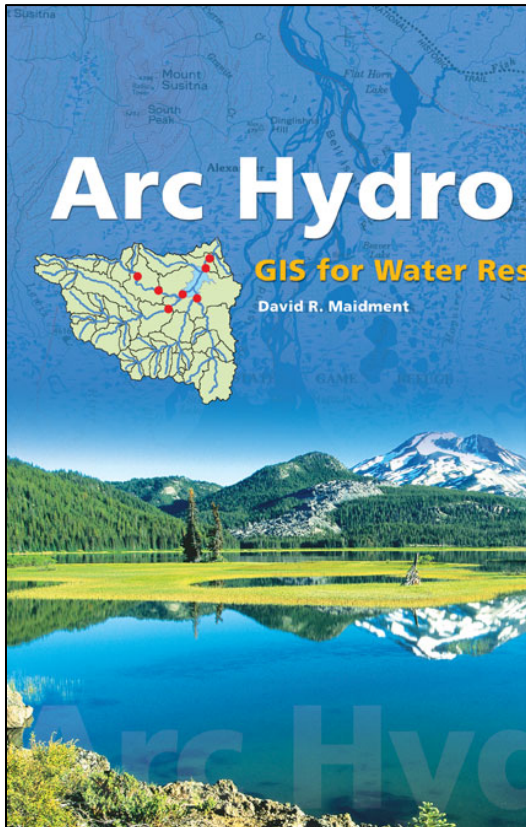
Sampling Feature

1	Stream Gage
---	-------------

Feature of Interest

1	Cross Section	Point
2	Stream	Line
3	Catchment	Polygon

Linking Data to the Geo-Environment



Arc Hydro
GIS for Water Resources
David R. Maidment

NHDPlus Version 2

Horizon Systems → NHDPlus Version 2

Horizon System Corporation

« NHDPlus Version 2 »

NHDPlusV2 consists of the following components:

- Greatly improved 1:100K National Hydrography Dataset (NHD)
- Greatly improved 30 meter National Elevation Dataset (NED)
- Nationally complete Watershed Boundary Dataset (WBD)
- A set of value added attributes to enhance stream network navigation, analysis and display
- An elevation-based catchment for each flowline in the stream network
- Catchment characteristics
- Headwater node areas
- Cumulative drainage area characteristics
- Flow direction, flow accumulation and elevation grids
- Flowline min/max elevations and slopes
- Flow volume & velocity estimates for each flowline in the stream network
- Catchment attributes and network accumulated attributes
- Various grids from the hydro-enforcement process including the hydro-enforced DEM.

Additional Information

NHDPlusV2 attribute and vector data are distributed by hydrologic regions (HUC2) or, in the case of regions 10 and 03, parts of hydrologic regions. These regions or parts of regions are called Watershed Processing Units (VPUs). The NHDPlusV2 raster components are distributed by sub-parts of VPUs called Raster Processing Units (RPUs).

The NHDPlusV2 User Guide is available [here](#).

Important information is also contained in the release notes that are posted with each VPU.

A national shapefile of the VPU and RPU boundaries will be available soon.

All available data is served from the Horizon Systems FTP site and may be freely **downloaded**.

NHDPlusV2 Pre-Release Data vs Public-Release Data

The first public release of NHDPlusV2 is in data model version 2.1. All of the names of the NHDPlusV2 public-release distribution files begin with "NHDPlusV21_". In parts of the country, the release version of NHDPlusV2 was provided to several collaborators that helped fund the NHDPlus production. These collaborators worked closely with the NHDPlus team to ensure that (1) the release data was appropriate for their use and (2) an eventual migration to the public version, NHDPlusV2 data model version 2.1, was possible. The pre-release version of the data is in data model version 2.0 and the names of the pre-release distribution files begin with "NHDPlusV2_".

Please note that the NHDPlus Team will be unable to provide support for the pre-release data, except to the collaborators referenced above. Our goal is to standardize future use on the public release version, NHDPlusV2 data model version 2.1.

OGC 11-039r2

Open Geospatial Consortium

Approval Date: 2012-03-23
Publication date: 2012-04-06

External identifier of this OGC® document: <http://www.opengis.net/doc/DP/hy-features>

Reference number of this OpenGIS® document: OGC 11-039r2

Category: OGC® Discussion Paper

Editors: Rob Atkinson, Irina Dornblut

HY_Features: a Common Hydrologic Feature Model Discussion Paper

Copyright notice

Copyright © 2012 Open Geospatial Consortium
To obtain additional rights of use, visit <http://www.opengis.org/legal/>.

Warning

This document is not an OGC Standard. This is an OGC Discussion Paper and is therefore not an official position of the OGC membership. The document is distributed for review and comment. It is subject to change without notice and may not be referred to as an OGC Standard. Further, an OGC Discussion Paper should not be referenced as required or mandatory technology in procurements.

Recipients of this document are invited to submit, with their comments, notification of any relevant patent rights of which they are aware and to provide supporting documentation.

Document type: OpenGIS® Discussion Paper
Document subtype: Encoding
Document stage: Approved for public release
Document language: English

Much of this work has been done already!

Capturing the Knowledge Content of Data

- What do the data mean?
- How have they been used?
- What conclusions have been drawn?
- What are appropriate uses?

Linking data to papers is great, but only a small fraction of data end up in papers!

Capturing the Knowledge Content of Data

Data Annotation

Comments



Amber Jones 9 minutes ago

These data characterize diurnal fluctuations in temperature in the Little Bear River at its terminus just upstream of Cutler Reservoir. At its lower end, the hydrology of the Little Bear River has been highly modified, with streamflows dominated by releases from Hyrum Reservoir and with several agricultural diversions and return flows. These temperature observations reflect a high degree of human modification within what was typically a spring snowmelt dominated system.

[Link](#) | [Reply](#)

Liked by 0 +1



Tony Melcher 6 minutes ago

I used these data in a fish habitat assessment for the Little Bear River. They demonstrate the diurnal variability in water temperature, and I was able to use them to assess potential acute and chronic effects of elevated stream temperatures on salmonid fish species.

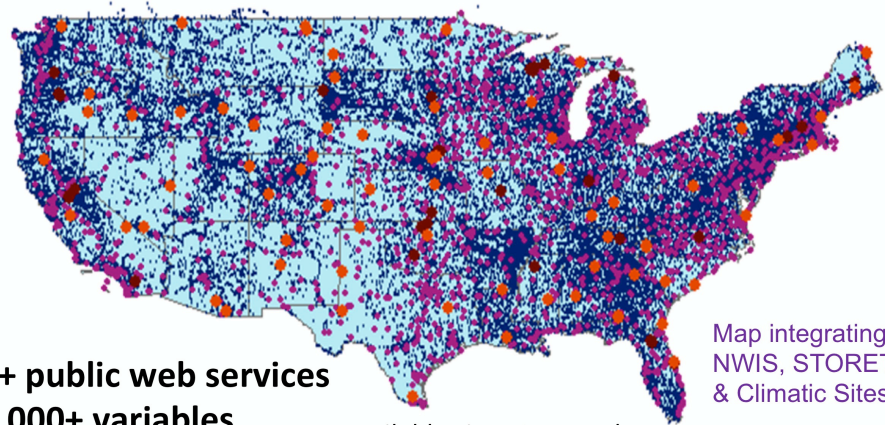
[Link](#) | [Reply](#)

Liked by 0 +1

Integration with Government Agency Data Repositories

- Many models require data at spatial and temporal scales not found in research data infrastructures
- For some science domains – there is far more data in agency repositories

HIS Central HydroCatalog Content



65+ public web services

13,000+ variables

1.96+ million sites

23.3 million observation time series

Referencing 5+ billion data values

*Available via HIS Central
discovery services*

Available via GetValues requests

Map integrating
NWIS, STORET,
& Climatic Sites

Metadata for most services are harvested weekly

Summary

- Speeding innovations in synthesis and modeling in the geosciences will require:
 - New CI techniques and tools
 - A work force capable of developing and using these tools
- Enhancing the discoverability, accessibility, and use of data in research repositories for modeling:
 - Common information modeling – standards promote interoperability
 - Linking data to the geo-environment for spatial context
 - Capturing the knowledge content of data

Questions?

Jeffery S. Horsburgh

Civil and Environmental Engineering

Utah State University

jeff.horsburgh@usu.edu