

Data Federation via Metadata

Bill Michener
University of New Mexico

EUDAT 2nd Conference
Rome, Italy
October 29, 2013

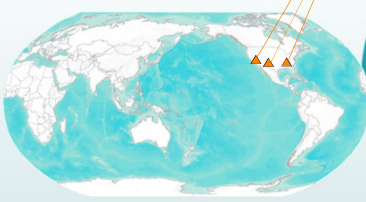


COLLEGE of UNIVERSITY LIBRARIES & LEARNING SCIENCES




Technology

Three major components for a flexible, scalable, sustainable network

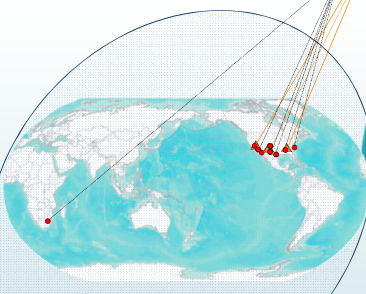


- Coordinating Nodes**
 - retain complete metadata catalog
 - indexing for search
 - network-wide services
 - ensure content availability (preservation)
 - replication services




Technology

Three major components for a flexible, scalable, sustainable network



- Coordinating Nodes**
- Member Nodes**
 - diverse institutions
 - serve local community
 - provide resources for managing their data
 - retain copies of data



Technology

Three major components for a flexible, scalable, sustainable network

Coordinating Nodes

Member Nodes

- diverse institutions
- serve local community
- provide resources for managing their data
- retain copies of data

Technology

Three major components for a flexible, scalable, sustainable network

Coordinating Nodes

Member Nodes

Investigator Toolkit

- command line interface
- python
- java
- GNE Drive
- DMP Tool
- GNE Mercury
- DATAUP
- Kepler
- GNE R

Preserve Data and Metadata

- Metadata mirrored at Coordinating Nodes
- Data replicated between Member Nodes
- CNs manage copies
- Checksums recorded and verified
- Promote quality metadata

Metadata Quality

- Data discovery
 - Title, abstract, investigators, etc.
 - Data extent (over space, time, and taxa)
 - Associated research papers
- Data interpretation and integration
 - Methodological details allow re-use
 - Semantics allow improved integration
- Analysis and modeling
 - Automate mechanical processing
 - Allows focus on science issues not plumbing



Metadata Formats

- Extract and index common fields from metadata standards
 - Ecological Metadata Language (EML)
 - FGDC Biological Data Profile (BDP)
 - ISO 19115 Geospatial Metadata
 - Dublin Core
 - Darwin Core
 - METS
- Extensible to include many more
 - DIF, NexML, WaterML, CF, NcML, ESML, DDI, MIENS, ...



Discover: MN Implications

- Supported science metadata formats
- Accuracy
- Extent
- Data package descriptions



Traditional Term Matching

**Poor Precision
Poor Recall**

Literal text matches miss hypernyms and hyponyms

query → term matching (TF-IDF) → ranked list of metadata records matching query

Can be weighted based upon how often a term appears in a specific document compared to how often it appears in an entire corpus (TF-IDF: Term frequency-inverse document frequency)

Adding Formal Ontologies

**Improved Precision
Improved Recall**

Query expansion based on formally-defined relationships between concepts

formal ontologies/
controlled vocabularies

query → term matching (TF-IDF) → ranked list of metadata records matching query

High-level and domain-specific ontologies specify formal and explicit representation of relationships between concepts (e.g., parent-child-sibling relationships; synonyms, partonyms, hypernyms)

Challenge: Linking data to ontology

OBOE Salmon Ontology

```

    graph TD
      SteelheadPopulationSample -- is-a --> AliveWildSmoltSteelheadPopulationSample
      SteelheadPopulationSample -- has-characteristic --> Count
      Count -- uses-standard --> Number
      Count -- has-precision --> One[1]
  
```

OBOE Semantic Annotation

```

    graph TD
      Observation1[Observation] -- has-measurement --> Measurement1[Measurement]
      Observation2[Observation] -- has-measurement --> Measurement2[Measurement]
      Measurement1 -- of-characteristic --> Count
      Measurement2 -- of-characteristic --> Count
  
```

Structural Metadata

```

    <attribute id="att.4">
    <attributeName>
    live_stlhd_smolt
    </attributeName>
    </attribute>

    <attribute id="att.5">
    <attributeName>
    live_count
    </attributeName>
    </attribute>
  
```

Data

date	live_stlhd_smolt	dead_stlhd_smolt	date	species	live_count	dead_count
2/15/09	27	2	2/15/09	STHD	45	0
2/16/09	34	0	2/16/09	CCHO	23	1

