

Life cycle semantics in Earth and space sciences – what's worked (and not) and where are we... a *decade* in...

Ontolog Forum (Earthcube) Sep. 6, 2012

Peter Fox (RPI) <u>pfox@cs.rpi.edu</u> Tetherless World Constellation







Scientists should be able to access a global, distributed knowledge base of scientific data that:
appears to be integrated
appears to be locally available

But... data is obtained by multiple means (instruments, models, analysis) using various protocols, in differing vocabularies, using (sometimes unstated) assumptions, with inconsistent (or non-existent) meta-data. It may be inconsistent, incomplete, evolving, and distributed. And, it is almost always created in a manner to facilitate its generation not its use.

And... there exist(ed) significant levels of semantic heterogeneity, large-scale data, complex data types, legacy systems, inflexible and unsustainable implementation technology... 2



Origins ...

- In 2000-2001 the need for capturing and preserving knowledge in science data became very clear but the barriers were high
- In 2004 we started a virtual observatory project based on semantic technologies
- Use case driven in solar and solar-terrestrial physics with an emphasis on instrument-based measurements and real data pipelines; we needed implementations
- We knew we also needed integration and provenance (but that came later)
- We aimed to push semantics into our systems to build new 'prototypes' but we 'failed' ;-)



In 2004

- 2004 OWL was a W3 recommendation!!
- Protégé 2.x and the Protégé-Java-OWL API
- SWOOP was a viable editor
- Jena and the Jena API were in good shape
- Pellet worked
- SPARQL was still a twinkle in the RDF working group's eye
- Semantics were still the realm of computer scientists – luckily we had one of the best

Design and Development

- We made a conscious decision only to develop ontologies that were required to answer specific use cases
- We made a conscious effort to use whatever ontologies were available**
- We were pretty sure that rules would be needed
- We ignored query



Use Case example

- Plot the neutral temperature from the Millstone-Hill Fabry Perot, operating in the non-vertical mode during January 2000 as a time series.
- Plot the neutral temperature from the Millstone-Hill Fabry Perot, operating in the non-vertical mode during January 2000 as a time series.
- Objects:
 - Neutral temperature is a (temperature is a) parameter
 - Millstone Hill is a (ground-based observatory is a) observatory
 - Fabry-Perot is a interferometer is a optical instrument is a instrument
 - Non-vertical mode is a instrument operating mode
 - January 2000 is a date-time range
 - Time is a independent variable/ coordinate
 - Time series is a data plot is a data product

Knowledge representation

Statements as triples: {subject-predicate-object}

interferometer is-a optical instrument

Fabry-Perot is-a interferometer

Optical instrument has focal length

Optical instrument is-a instrument

Instrument has instrument operating mode

Instrument has measured parameter

Instrument operating mode has measured parameter

NeutralTemperature is-a temperature

Temperature is-a parameter

- A query*: select all optical instruments which have operating mode vertical
- An inference: infer operating modes for a Fabry-Perot Interferometer which measures neutral temperature
- ISWC paper award 2006, IAAI best paper (2007), Fox et al. 2009 in Computers and Geosciences.

Semantic Web Benefits

- Unified/ abstracted query workflow: Parameters, Instruments, Date-Time
- Decreased input requirements for query: in one case reducing the number of selections from eight to three
- Generates only syntactically correct queries: which was not always insurable in previous implementations without semantics
- Semantic query support: by using background ontologies and a reasoner, our application has the opportunity to only expose coherent query (portal and services)
- Semantic integration: in the past users had to remember (and maintain codes) to account for numerous different ways to combine and plot the data whereas now semantic mediation provides the level of sensible data integration required, and exposed as smart web services
 - understanding of coordinate systems, relationships, data synthesis, transformations.
 - returns independent variables and related parameters
- A broader range of potential users (PhD scientists, students, professional research associates and those from outside the fields)

Semantics - Modern informatics enables a new scale-free** framework approach

Semantic Web Methodology & Technology Development Process

- Use cases
- Stakeholders
- Distributed authority
- Access control
- Ontologies
- Maintaining Identity

- Establish and improve a well-defined methodology vision for semantic technology based on application development
- Leverage controlled vocabularies, etc.





http://www.w3.org/2003/Talks/1023-iswc-tbl/slide26-0.html, http://flickr.com/photos/pshab/291147522/

Developing ontologies (c. 2005)

- Use cases and small team (7-8; 2-3 domain/ data experts, 2 knowledge experts, 1 software engineer, 1 facilitator, 1 scribe)
- Identify classes and minimal properties (leverage controlled vocab.)
 - Start with narrower terms, generalize when needed or possible
 - Adopt a suitable conceptual decomposition (e.g. SWEET)
 - Import modules when concepts are orthogonal
 - Add service classes and properties where needed
- Review, vet, publish
- Only code them (in RDF or OWL) when needed (CMAP, ...)
- Ontologies: small and modular



Semantics between 2004 and 2009

- Ontologies were needed for data integration and provenance and mediation for data mining
- Protégé 3.x and then 4.0 came out
- SWOOP development was interrupted
- Cmap added OWL predicate support*
- SPARQL became a recommendation
- Triple stores exploded in use and capability
- Linked Open Data started to take off
- Pellet 2.0 came out
- We invaded OWLED 2006, 2007, 2009, (2010) Tetherless World Constellation



Working with knowledge

Expressivity

Implementability

Maintainability/ Extensibility





Or it may be this ...



Rule execution





Semantics between 2009 and 2012

- Semantic data framework (SeSF)
- Substantial knowledge provenance work
- Data quality, uncertainty and bias representations and applications (oh, these are in production at NASA)
- Multi-sensor advisor:







RuleSet Development

[DiffNEQCT:

->

(?s rdf:type gio:RequestedService),

(?s gio:input ?a),

(?a rdf:type gio:DataSelection),

(?s gio:input ?b),

(?b rdf:type gio:DataSelection),

(?a gio:sourceDataset ?a.ds),

(?b gio:sourceDataset ?b.ds),

(?a.ds gio:fromDeployment ?a.dply),

(?b.ds gio:fromDeployment ?b.dply),

(?a.dply rdf:type gio:SunSynchronousOrbitalDeployment), (?b.dply rdf:type gio:SunSynchronousOrbitalDeployment), (?a.dply gio:hasNominalEquatorialCrossingTime ?a.neqct), (?b.dply gio:hasNominalEquatorialCrossingTime ?b.neqct), notEqual(?a.neqct, ?b.neqct)

(?s gio:issueAdvisory giodata:DifferentNEQCTAdvisory)]



Multi-sensor Data Synergy Advisor (NASA), Leptoukh, Lynnes, Zednik, et al.

Semantic Advisor Architecture



Multi-sensor Data Svnerov Advisor (NASA). Leptoukh. Lvnnes. Zednik. et al.



Semantics between 2009 and 2012

- Semantic data framework (SeSF)
- Substantial knowledge provenance work
- Data quality, uncertainty and bias representations and applications (oh, these are in production at NASA)
- Multi-sensor advisor
- Applications:

 Sea Ice, Carbon Observatory, Integrated Ecosystem Assessments, globalchange.gov, ocean.data.gov, energy.data.gov 21 Respect and Mediation ... how

O O O O O O O O O O O O O O O O O O O	Biological and Chemical (er/mapsdev-ol/index.php	Oceanography Data Collection		¢ Qr Bing	O	
BC - DM MapServer Geospat	al Interface			Contact Help NS	F Acknowledgment 😽	
Search	Results					
▼ Programs	Available datasets			Mapped datasets		
	Group by: dataset		v	S Remove all		
	Dataset	Deployment				
(2358) U.S. GLOBal ocean ECosystems dynamics (903) U.S. Joint Global Ocean Flux Study	Dataset: 2dmodel_m2only					
(881) National Marine Fisheries Service / Northeast (503) Ocean Carbon & Biogeochemistry	C 2dmodel_m2only lab_UNH-2dmodel					
(449) Census of Marine Life (213) Iron Synthesis	□ Dataset: 3-D Basin-Scale Ecosystem Model of the North Atlantic			No datasets have been mapp	ped.	
(E1) United States Surface Ocean Lower Atmember	3-D Basin-Scale Ecosystem Model of	the North Atlantic USJGOFS_SM	/P	Click the plus icon next to a datase	t to begin.	
(51) United States Surface Ocean Lower Atmospher	atasati 2-D Ecosystem Model of the	Poss Con				
▼ SeaVox Categories	3-D Ecosystem Model of the Ross Se	a USJGOFS SM	/P			
(1192) discrete water samplers	ataset: 3d_model_simulations		_			
(1056) plankton nets (761) CTD profilers	A Page 5/9 of 579 ▶ ▶ ■ C	5781 - 5788 of 5	788			
(591) Sea-Bird SBE 19 SEACAT CTD (550) Sea-Bird SBE 911 CTD	V le deployments Map					
(428) Sea-Bird SBE 911plus CTD	Clear selections	🕂 Pan 🕕 Query map 🥠 Clear query			Map options -	
(396) thermosalinographs	A16N_33RO20030		A DEREN	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	1 Same	
	AB_63_1	ASIA	1- Ar-		ASIA	
▼ Instruments	AB_63_2	EUROPE		NORTH A	EUROPE 45900'N	
	AB 63 4A	May A A A	1 00	AMERICA	alay 1	
(1157) Niskin Bottle (591) CTD Sea-Bird SEACAT 19	AB_63_A		19 17			
(560) CTD Sea-Bird 911	AB_64_5		(X) \$			
(468) Bongo Net (460) CTD Neil Brown Mark V	AB_64_6			N AMERICA		
(428) CTD Sea-Bird SBE 911plus (408) CTD Neil Brown Mark III	AB_64_7	AUSTRALIA				
(395) Thermosalinograph	AB_04_0		2		45:00'S	
	AESOPS_Array			VI KES	The state of	
▶ People	All-119-4	SQUTHERN OCEAN E 45°00'E 90°00'E 135°00'E	180907 135°001	w. 90°00 w 2 2°00 w 00°00	68.5282° S 37.3125° W 45°00'E 90°00'E	



	Results				
rograms	Available datasets		Mapped datasets		
rojecte	Group by: deployment	~		😣 Remov	
	Dataset	Deployment			
spie	Deployment: AE-X1	103			
	CTD Profiles	AE-X1103			
nt Categories	Deployment: ATI-119-4		No datasets have be	en mapped.	
		All-119-4	Click the plus icon next to a	dataset to begin.	
(437) Sea-Bird SBE 911plus CTD		All-119-4			
(184) fluorometers		All-119-4			
(1) Turner Designs SCUFA II Submersible					
(20) WETLabs ECO EL fluorometer	Deployment: AII-11	9-5			
(20) WEILEBS LEO I'E Indofonitetei	Page 1 of	19 🕨 🕨 🦧 1 - 25 of 184			
fluorometer	Visible deployments	Man			
(8) WETLabs ECO-FLNTU combined	Clear selections			Man antia	
(40) current meters	AE-X1103	Cicar query			
	All-119-4	A			
ments	All-119-5	ASIA			
ter Categories	AL9701		AMERICA	EU #45	
·	AL9705			THE TLANTIC	
ers	AL9707		DCEAN	OCEAN AFRICA	
\$	AL9806		PACIFIC	00	
Platforms	AL9808	191 103	OCEAN	AMERICA	
	AL9901	INDIAN AUSTRALIA			
	AL9904	OCEAN	D-	2 4	
	AL9906			SA I	
	T AT44 47	The second s			

Core and Framework Semantics -Multi-tiered interoperability







Summary

- In 2004 we set out to build a prototype and ended up with a production semantic data framework
 - Languages and tools served us well
- Even with modest expressivity we challenged the tools of the time and made many compromises
- All along the way, we evaluated our ontology developments and implementations to gauge the benefits of semantics
- Maintainability, esp. modularization is driving new expressivity needs
- Xinformatics and a repeatable methodology is the key (information models) - we continue to need to bridge computer science and application communities ("It's the language stupid", i.e. semantics)



Contact

- pfox@cs.rpi.edu
- http://tw.rpi.edu
- @taswegian
- See also wiki.esipfed.org (Semantic Web Cluster)