

Knowledge Provenance

CMSP Vocabulary Workshop Dec 1-3, 2010 (Woods Hole, MA)

Peter Fox (RPI) pfox@cs.rpi.edu

Tetherless World Constellation – http://tw.rpi.edu





Provenance

- Origin or source from which something comes, intention for use, who/what generated for, manner of manufacture, history of subsequent owners, sense of place and time of manufacture, production or discovery, documented in detail sufficient to allow reproducibility
- Or ... provenance are the types and instances of metadata in a particular multi-faceted context
- Knowledge provenance; enriched with semantics (especially the relations between concepts previously isolated, and retaining context) and semantically-aware tools



Problem definition

- Data is coming in faster, in greater volumes and outstripping our ability to perform adequate quality control
- Data is being used in new ways and we frequently do not have sufficient information on what happened to the data along the processing stages to determine if it is suitable for a use we did not envision
- We often fail to capture, represent and propagate manually generated information that need to go with the data flows
- Each time we develop a new instrument, we develop a new data ingest procedure and collect different metadata and organize it differently. It is then hard to use with previous projects
- The task of event determination and feature classification is onerous and we don't do it until after we get the data

Data has Lots of Audiences





Fact: Scientific data services are increasing in usage and scope, and with these increases comes growing need for access to provenance information.

Our Research Goal: design and implement an extensible provenance solution that is deployed at the science data ingest/ product generation time.

Provenance Infrastructure Goal: to design a reusable, interoperable provenance infrastructure.

Outcome: implemented provenance solution in one science setting AND operational specification for other scientific data applications – both achieved.



Explanation Justification Verifiability



Proof

Trust



But back to reality Fragmentation Disconnection - Internal/Ext. Encapsulation

... all are bad for ... transparency



000

MLSO Analysis Center Homepage



Welcome to the Mauna Loa Solar Observatory (MLSO) Website. The MLSO, operated by the High Altitude Observatory in Boulder Colorado, houses several instruments designed to observe the sun at many different wavelengths. The MLSO instruments provide observations needed to understand the sun's continuous release of plasma and energy into interplanetary space.

11.

Advanced Coronal Observing System. A suite of instruments designed to observe the solar atmosphere at a variety of heights. Includes Chromospheric Helium Imaging Photometer (CHIP, 1083.0nm), H–alpha prominence and solar disk monitor (PICS, 656.2nm), and the Mk4 K–coronameter, which observes the white light K–corona from 1.12–2.79 solar radii. ECHO Experiment for Coordinated Helioseismic Observations. A network of two instruments which observe solar oscillations as seen in the radial velocity of the solar surface. PSPT PSPT PSPT PSPT

MLSO Home	ACOS Home	ACOS Data	ECHO Data	PSPT Data	Mauna Loa Webcams	Related Sites	Contact Us	Eclipses	Publications	About MLSO





- Typical science data processing pipelines
- Distributed
- Some metadata in silos
- Much metadata lost
- Many human-in-loop decisions, events
- No metadata infrastructure for any user





The ACOS case for Provenance

- Provenance metadata currently not propagated with or linked to the data products
 - Processing metadata
 - Origin (observation) metadata
- Data products are the result of "black box" systems
 - Most users do not know what calibrations, transformations, and QA processing have been applied to the data product



Processing

10



- What were the cloud cover and seeing conditions during the observation period of this image?
- What calibrations have been applied to this image?
- Why does this image look bad?

Why does this image look bad?



Provenance and Domain concepts in the use cases





Multi-domain Knowledge Base



Proof Markup Language (PML)





PML NodeSet









Open Provenance Model

- Agents
 - Catalyst and controlling entity of a process
- Processes
 - Action or Series of actions performed resulting in new artifacts
- Artifacts
 - Immutable piece of state
- Roles
 - Non-semantic flat tags used to provide context in relations



Concept Alignment (OPM)







Multi-Model Individuals

Individuals use OWL's support of multiple inheritance to provide mediation between the models







PML NodeSet using Multimodel individuals



Knowledge Base with Provenance and Domain Models in Alignment



Alignment via Ontology Constructs

- Use ontology constructs to map a relationship between concepts in different domains
- Can be defined in a separate ontology than the models being mapped
- Does not require a change to the source models!
- OWL
 - owl:equivalentClass
 - owl:equivalentProperty
 - owl:sameAs

RDFS

- rdfs:subClassOf
- Rdfs:subPropertyOf



Direct Alignment using Rules*

- Rules provide conditional logic on semantic constructs outside application logic
- Rules can be updated or tweaked without requiring an application update.
- Easily shared and managed
- Provides for more complex mapping than ontology constructs

ex:Instrument(?x)

➔ pmlp:Sensor(?x)

pmlp:Information(?x) ^

pmlp:hasURL(?x,?url) ^

swrlb:endsWith(?url, ".hsh.fts ")

→ Ex:CHIPIntensityImage(?x)

*Many rule systems exist, this slide uses the Semantic Web Rule Language (SWRL)

Querying/Interrogating the Knowledge Base

• Back to one use case:

What *calibrations* have been applied to this *image*?

- We construct a query returns any individuals with type Calibration used as the InferenceRule in the justification from any artifact the current artifact was derived from.
- We assume that any calibration applied to an artifact the current artifact was derived from can also be considered as 'applied' to the current artifact, and that the wasDerivedFrom property is transitive





Provenance Capture





Provenance aware faceted search





Interoperability with Provenance Tools





Multi-sensor Synergy Data Advisor (MDSA)

- Based on NASA sensor assets and multi-option processing -<u>http://giovanni.gsfc.nasa.gov</u>
- Dynamically generated lineage (XML)
- Want to advise for or against certain processing operations on certain data products (both internal and external provenance)

Giovanni Allows Scientists to Concentrate on the Science





Inter-comparison of data from multiple sensors

Data from multiple sources to be used together:

- ACE
- NPP and NPOESS
- Geo-Cape
- European and other countries' satellites
- Models

Harmonization:

- It is not sufficient just to have the data from different sensors and their provenances in one place
- Before comparing and fusing data, things need to be harmonized:
 - Data: format, grid, spatial and temporal resolution
 - Metadata: standard fields, units, scale, quality
 - Provenance: source, assumptions, algorithm, processing steps

Dangers of easy data access without proper assessment of joint data usage - It is easy to use data incorrectly

Why don't MODIS Terra and Aqua Aerosols agree?



Sensitivity Study: Daily AOD MODIS Terra vs. MISR Terra

Collecting and Delivering Data Provenance

Where to find the knowledge about data and data

processing?

- It is scattered in scientific papers, the actual code, unwritten assumptions, folklore, etc.
- Assess sensitivity of the results to variations in processing algorithms/steps...
- Work closely with scientists to guarantee science quality

How to deliver provenance?

- Deliver to users together with the data
- Present to users in a convenient, easy-to-read fashion
- Provide recommendations for different data usage (applications vs. climate studies)

Use Case Processing PML

• Proof Markup Language graph for Giovanni processing provenance

 Capability to access PML using IW Browser or Probe-It provenance visualization tools.

Accessing PML using Probe-It provenance visualization tools

IWBrowse

However...

- Normal people cannot look at these presentations and get answers to their questions, and often they don't care!
- Or, if they do, then they want to see a semantic provenance difference - no they don't really know to ask for this but it is the step to establishing trust in these systems
- Provenance presentation is still an open challenge... 43

So ... Semantic Advisor

Semantic Advisor

Provides caveats for intercomparison based on user selections

- Parameter, Dataset
- Satellite, Orbit (derived from Dataset)

Service-based Interaction

- Giovanni sends XML with user input
- RPI returns XML with comparison info
- Giovanni renders comparison as table

Features

Leveraging Ontology and Rulesets

- Developed ontology to represent advisories and corresponding semantic rule set to assert advisories when certain conditions hold in the knowledge model.
- Developing web service to ingest user input XML, translate into semantic knowledge based on ontology vocabularies, perform reasoning and rule-based inference on the knowledge, and generate a response XML to return to the Giovanni service.

RuleSet Development

Rulesets:

- Rulesets are used to understand the criteria that make up a fitness for purpose or important factors for assessing the concept of a 'set of rules
- Individual rules need to be captured so that the significance of different parameter combinations can be explained / presented to an end user
- Rulesets are needed to indicate to a computer what information important and needs to be propagated in/ to the processing results.

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)]

Semantic Advisor

Your Selected Ontions:

Semantic Advisor

Spatial Area: Parameters:

Temporal Range

Longitude (-30, 150), Latitude (-10,60) A: MYD08_D3.005 Aerosol Optical Depth at 550 nm B: MOD08_D3.005 Aerosol Optical Depth at 550 nm Begin Date: Jan 01 2008

	Tomporar Rango.	Dogin Date. Buildi 20			
Parameter Name :	Aerosol Optical Depth at 550 nm	Aerosol Optical Depth at 550 nm			
Dataset:	MYD08_D3.005	MOD08_D3.005	← Diff		
Data-Day definition	UTC (00:00-24:00Z)	UTC(00:00-24:00Z)	The same but		
Temporal resolution	Daily	Daily			
Spatial resolution	1x1 degree	1x1 degree			
Sensor:	MODIS	MODIS			
Platform:	Aqua	Terra	← Diff		
EQCT	13:30	10:30	← Diff		
Day Time Node	Ascending	Descending	← Diff		
Pre-Giovanni Processes :	ATBD-MOD-30	ATBD-MOD-30			
Giovanni Processes:	Spatial subset Time average	Spatial subset Time average			
		10 6 2 30 60 90 120 150	a 2 20 -60 -30 (
MODIS Terra vs. MODIS Aqua A	OD Correlation	Included Overpass time Difference			
Continue process to d	isplay image	Return to selection page			

Final Remarks / Discussion

- Integrated Knowledge Provenance and Domain Knowledge Base key to our Use Cases
- PML supports
 - Multi-model individuals (by way of OWL)
 - Causality graphs and justifications
 - Processing history and intent
- Multi-model individuals interoperable with generic PML tools
 - PML is somewhat hard to generate unless you are an expert
 - Tools are needed around this
- Ongoing Steps
 - Investigation of using Rules to infer domain relations from provenance store (and vice versa)
 - Further development on Semantic Faceted Browse
 - Further design on visualization of provenance

Acknowledgements

- RPI + Inference Web team
 - Stephan Zednik
 - Deborah McGuinness
 - Patrick West
 - James Michaelis
- NASA/GSFC
 - Gregory Lepkotuh
 - Chris Lynnes
- UTEP/CyberSHARE
 - Paulo Pinheiro da Silva
 - Nicholas del Rio

Links

- PML: http://inference-web.org/2007/primer/
- OPM: <u>http://openprovenance.org/</u>
- SWRL: <u>http://www.w3.org/Submission/SWRL/</u>
- Inference Web http://inference-web.org
- Probe-It! <u>http://trust.utep.edu/probe-it/</u>
- SPCDIS <u>http://tw.rpi.edu/portal/SPCDIS</u>
- MDSA <u>http://tw.rpi.edu/portal/MDSA</u>
- Many others...
- Contacts:
 - pfox@cs.rpi.edu
 - <u>zednis@rpi.edu</u>
 - westp@rpi.edu

Syntactic Temporal Constraint

Semantic Temporal Constraint

