An open-source information model and associated cyberinfrastructure for effective environmental management and analytics

Roelof Versteeg¹, Rebecca Rubinstein¹, Reza Soltanian² and Doug Johnson¹

1: Subsurface Insights 2: University of Cincinnati



2023 Global Summit on Environmental Remediation @REMPLEX



Acknowledgement

- The work presented here was funded by the DOE SBIR program under multiple SBIR awards (DE-SC-0009732, DE-SC-0018447 and DE-SC0019621)
- Examples shown include work performed by Subsurface Insights funded under the LBNL Watershed SFA and the SRNL ALTEMIS program, as well as from work done with collaborators from ILVO

Outline

- <u>Subsurface Insights introduction</u>
- Setting the stage: why is data management and analytics to understand contaminated site behavior still challenging
- Solution attributes
- Solution implementation: ODMX
- Examples
- Summary

Subsurface Insights introduction

 Geoscience company specializing in near surface process monitoring



- Main products include
 - Autonomous electrical resistivity systems and associated cyberinfrastructure for continuous subsurface monitoring [Contact me after my talk if interested]
 - End to end cloud-based software framework for geoscience data ingestion, management, analytics and reporting [THIS TALK]



End to end cloud-based software framework for geoscience data ingestion, management, analytics and reporting: includes data acquisition, transmission, data management, analytics and reporting



Outline

- Subsurface Insights introduction
- <u>Setting the stage: why is data management and analytics to understand</u> <u>contaminated site behavior still challenging</u>
- Solution attributes
- Solution implementation: ODMX
- Examples
- Summary

Truism: Contaminated site behavior is controlled by physical, chemical and biological natural and anthropogenic processes at multiple spatial and temporal scales





<u>Reality</u>

- different units and parameter names for same parameters (airtemp_c, tempair_F)
- timestamps in different time zones and formats
- different data formats and data density
- different meaning and "value" of different datasets

Site data is heterogeneous in type, provenance and meaning ("how data relates to processes")

- **Process Drivers**: precipitation, solar radiation, regional groundwater gradients, contaminant releases, remedial efforts
- **Process Context:** Elevation, Soil type, Geographic location
- **Process manifestation:** Vegetation, soil microbiology, soil moisture, contaminant concentrations

Large provider with dedicated infrastructure generating domain specific data	Vendor hosted platforms for sensor data	Laboratory Information Management Systems (LIMS)	In house data collection systems
USGS: NWIS	Meter	In house	Data loggers with manual download
NASA: Landsat	In Situ	EMSL	
ESA: Sentinel	Campbell	JGI	Remediation systems (SCADA)
DOE ARM	Vaisala	Commercial	Sample
USDA: SNOTEL	SSI ERT systems		systems

Example: Data heterogeneity in type, provenance and meaning

The data management challenge:

integrate the multiple inter-related datasets that provide complementary views of the same system so that we can use this data effectively for process understanding and site operation

Note: data management is an enabler

Three main requirements

- 1. Data needs to be normalized
 - Same parameter collected by two groups can have different names and/or different units
 - Data can be reported in different coordinate systems
 - Data can be reported in different time zones
 - Data organization and structure differs between different providers
 - → Normalize parameter names, units, coordinate systems and timestamps
- 2. Data needs to be fused
- 3. Data needs to be findable, accessible and usable by analytics codes

What is data fusion* and why do we need it

- **Data fusion**: an approach for combining and interconnecting multi modal data by providing relations between different data sets
- For some data existing relations are sufficient
- For some data we need to define additional relations

* In this context - there are other definitions of data fusion



Relations can be defined at the data management and at the analytics level

Outline

- Subsurface Insights introduction
- Setting the stage: why is data management and analytics to understand contaminated site behavior still challenging
- Solution attributes
- Solution implementation: ODMX
- Examples
- Summary

Attribute 1 of solution: represent complex relationships and attributes, support heterogeneous data types and ways to fuse data



Attribute 2 of solution: effective ways to ingest and extract data and allow for data discovery

Pipelines for data harvesting, processing and ingestion

Database which can holds heterogeneous data with rich metadata

APIs for data discovery and retrieval

Attribute 3 of solution: easy to deploy and use, flexible and extensible

Robust software tools allowing users to write their own parsers Extensible (e.g. project specific Controlled Vocabularies)

Cloud deployable

Two part solution which has these attributes: ODMX



ODMX Information Model

ODMX Cyberinfrastructure

Based on powerful NSF funded ODM2 Information Model (Horsburgh, 2016) Developed under SBIR funding. Deployed and used by e.g. LBNL SFA and SRNL ALTEMIS

Propie, Organizations & Affiliations		=Observations = Actions + Results	ODMICore.Variables	aDatasets & Citations			
ODMOCore Persons	OP ODHKCare. Organizations	DOMOCare, Methods	VatableD	ODMECare.Processis	gLevels ODMECore.Datasets		
Arsonifi Varie Prosofit Varie Prosofit Varie Prosofit Varie Organization Organization Organization Organization Organization Organization Organization Organization		Harthold Helhold year Michold rate Historiatie Historiatie	VariableReme VariableReme VariableOefnitten VariableOefnitten VariableOefnitten SampledhedumCV B - GuaritskindCV B - nableSocrecUR B	Processing.evelD Processing.evelCode Definition Explanation	D Dataset00 Dataset0400 Dataset7ypaCU Dataset7ypaCU Dataset7this Dataset8htmat		
controles, Amilations	H	LIGHTLAN, OLD , MILLIONS	*	CODMERCE	re.DatasetsResults +		
Artission Artission Cogenization Istransportanisti (Contact Istransportanisti (Contact Istransportanisti Affiziazione Affiziazione Statefrownis Statefrownis Statefrownis Costry Personalité	COMOCine Arliandy	Amorio Amorio Amorosovieta Amorosovieta Benduat Time Utomat Benduat Time Utomat Benduat Time Utomat Benduat Time Utomat Benduat Time Utomat Benduat Time Utomat	CONTRACTOR AND A CONTRA	Cite BridgeD Dutatett Result0 Dutatett Result0 Dutatett Result0 Dutatett Result0 Dutatett Result0 Dutatett	Struce name of attacet (itations HID HID HID HID HID HID HID HID HID HID		
HUTE C		Related coheletionD D	ValueCount		Publisher		
and and a service of		W. Contraction of the second s	Netherarde		CtationLink II		
OCOECore.SamplingFeatures	ODMCS ampling Features. Spatial Offse	sRoutes & Values	+	± Ψ			
anzingFeatureD	SpatialOffsetD		+	T 00M	Diffesults.ProfileResults		
lamping/seture(3)6 - imping/seture(3)et/ lamping/seture(3)et lamping/seture(3)et lamping/seture(3)et lamping/seture(3)et lamping/seture(3)et lamping/seture(3)et lamping/seture(3)et lamping/seture(3)et lambi	Spatial/HartgecV Offset(Vet0 Offset(Vet0 Offset2Vet0 O	Concessus Management/Results ResultD Resu	 Domanesults. Tensiseri ResultD Xuoation Viocation Viocation 22:eeat	Approximate A	0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1		
(PSG 0	Relationsh	× x	COMVRand's Lines	and and the second second	Company runs restorates		
Constanting of the second of t		COMMING SUBLIC Feasibility of the second memory for sublic second memory of the second m	ValueD ResultD Dataviale ValueDat Time/TCoThet DeteoCode/Time/TCoThet DeteoCode/Time/TCoThet DeteoCode/Time/TCoThet Time/agregation/Troinal Time/agregation/Troinal		VisikaB MesikB Ostavaba VisikaBashtme Visika		

Table: Actions

Describes actions (e.g., observation, sample collection, sample analysis, field visits, field activities, etc.)

Columns

РК	Name	Data Type	NULL	Auto	Default	Comment	FK
67	ActionID	Integer		~	nextval("ODMXCore"."Actions_ActionID_seq"'::regclass)	Primariy key identifier	
	ActionTypeCV	VarChar()				CV term describing the type of action (e.g., observation, sample collection, sample analysis)	-
	ActionName	VarChar()	~			A name for an Action or activity.	
	ActionDescription	VarChar()	~			Text that describes the action	
	MethodID	Integer				Foreign key identifier for the method used to complete the action	~
	BeginDateTime	DateTime				The date/time at which the action began	
	BeginDateTimeUTCOffset	Integer				The UTCOffset for the BeginDateTime	
	EndDateTime	DateTime	~			The date/time at which the action ended	
	EndDateTimeUTCOffset	Integer	~			The UTCOffset of for the EndDateTime	
	ActionFileLink	VarChar()	1			A URL or path to a file created by or used by the Action, such as instrument output or configuration.	

ODMX Information Model

- 196 tables
- Can represent complex relationships, hold data, trace provenance
- Flexible extensions
- Includes model for microbiological data
- Substantially enhanced and improved from ODM2

ODMX Cyberinfrastructure

Pipelines for data harvesting, processing and ingestion

Database building and population tools from JSONs APIs for data discovery and retrieval

Controlled vocabulary in JSON format (easily modifiable and extensible)

Library allowing for easy development

Key aspect of data model: sampling features

- Sampling Features defined in the Open Geospatial Consortium Observations, Measurement and Samples standard (Cox 2010).
- "geospatial or physical entities on which or at which observations are made".
- All data in ODMX is associated with a Sampling Feature.

Sampling feature example

- Site
- Instrumented borehole
- Sample (specimen) and sub-specimens

Data is associated with sampling features, and sampling features are associated with each other

Sub-specimens are a child of Specimens, which was collected at an Instrumented Borehole, which is contained in a Site

Sampling features attributes

- Sampling Features can be hierarchical
- Sampling feature relations allow us to add additional level of information which provides relations between data
- Relations expressed through extensible controlled vocabulary, e.g.
 - Is Contained in
 - Is Sub Specimen of
 - Is Child of
 - Is Derived from
 - Is Part of
 - Is Related to
 - Was collected At
- Relations between sampling features provide a natural way to fuse data

Examples LBNL SFA Site











 $\leftarrow \rightarrow$

Report

Models

in

Samples | Subsurface Insights

Location

ER-MCP5

ER-MDP1

ER-MDP2

ER-MDP3

ER-MMT1

ER-PHF0

ER-PLM1

ER-PLM11

ER-PLM12

Name ER-

ER-

Samples at ER-PHF0 on 5 4 2014

PHF0_2014_05_04_12_00_00_icpms_subspecimen

PHF0_2014_05_04_12_00_00_anion_subspecimen

EK-INY I



# of samples		Specimen	IS	search			
209		Name	Name				
4		ER-PHF0_2	ER-PHF0_2014_05_02_12_00_00_collected				
1		ER-PHF0_2	ER-PHF0_2014_05_03_12_00_00_collected				
1 124 2184 252 18 17		ER-PHF0_2	ER-PHF0_2014_05_04_12_00_00_collected				
		ER-PHF0_2	ER-PHF0_2014_05_05_12_00_00_collected				
		ER-PHF0_2	ER-PHF0 2014 05 06 12 00 00 collected				
		ER-PHF0_2	ER-PHF0 2014 05 07 12 00 00 collected				
		ER-PHF0_2	ER-PHF0 2014 05 08 12 00 00 collected				
		ER-PHF0 2	ER-PHF0 2014 05 09 12 00 00 collected				
		FR-PHF0 2	FR-PHE0 2014 05 10 12 00 00 collected				May 10 2014
	search	Sample ER- PHF0 2014 05 04 12 00 00 icpms subspecimen					
	Description	Analyte	Value		Std Dev		Units
nen	Subspecimen of ER- PHF0_2014_05_04_12_00_00_collected	Li	0.746999		0.05900559		ppb
	Subspecimen of ER- PHF0_2014_05_04_12_00_00_collected	В	3.389706		0.11921479		ppb
nen		Na	1378.756579		35.58669468	ppb	
		Mg	7217.630019292		140.380532850635		ppb
		AI	22.809252		.52407244		ppb
		Si	2087.004965		46.76930719		ppb
		К	477.832012		1.43885408		ppb

90% 🖒

ப

 \bigtriangledown

: Log out Roelof Versteeg

 \gg \equiv

 $\leftarrow \rightarrow$ @



☆

90%

Cross period | Subsurface Insights







ALTEMIS project (Savannah River)

Well Optimization



 $\leftarrow \rightarrow \bigcirc$

in M € Log out Log out Cersteeg

Multi-Location Visualization | Subsurface Insights





Integration with analytics codes

- Analytics codes can query ODMX through python API allows for data and relation discovery and retrieval
- ODMX and associated databases are used by Subsurface Insights analytics codes
- Example is the use by the ALTEMIS project of ODMX for machine learning using PyLEnM (Meray et al, 2022)
- Same capability used by electrical resistivity processing codes

Summary

- Effective environmental datamanagement requires
 - Normalizing multi-modal data
 - Fusing multi-modal data
 - Ensuring that data can be easily used by analytics codes
- ODMX provides a robust and extensible solution which meets these requirements

https://github.com/subsurfaceinsights/odmx

https://odmx.org

Questions?

Roelof Versteeg 603 443 2202 <u>Roelof.Versteeg@subsurfaceinsight.com</u> <u>www.subsurfaceinsights.com</u>

ODMX additional details:

- Github repo has example on how to build from scratch
- Includes automatically generated library which matches information model
 - All functions and attributes accessible through autocomplete
 - Includes type checking
 - Allows users to quickly write their own parsers

```
nport odmx.support.db as db
 mport odmx.data_model as odmx
                                                          (sampling_feature_uuid: str, sampling_feature_type_cv: str,
                                                          sampling_feature_code: str, sampling_feature_id: int | None = None,
con = db.connect( : Connection
                                                          sampling_feature_name: str | None = None, sampling_feature_description: str
   db host='localhost',
                                                          None = None, sampling_feature_geotype_cv: str | None = None,
                                                          feature_geometry: str | None = None, feature_geometry_wkt: str | None = None,
   db_name='odmx_example',
   db user='odmx',
                                                          elevation m: float | None = None, elevation datum cv: str | None = None,
                                                          latitude: float | None = None, longitude: float | None = None, epsg: str |
   db_pass='odmx'
                                                          None = None, PRIMARY_KEY: str = 'sampling_feature_id') -> None
with con.transaction():
                                                          sampling_feature_uuid
   new_sampling_feature_names = ['SF1', 'SF2', 'SF3']
   sampling_feature_objs = [] : list[Unknown]
                                                          Describes the sampling features on which observations are made.
   for name in new_sampling_feature_names:
       sampling_feature_objs.append(odmx.SamplingFeatures(name=name))
```