# Cross-Repository Data Integration using Ontology Design Patterns

Adila Alfa Krisnadhi

DaSe Lab for Data Semantics
Wright State University

# About DaSe Lab

- 2 Faculty members: Dr. Pascal Hitzler & Dr. Michelle Cheatham

- 5 full time PhD students (+ a few master's and part-time PhD students)

- Topics:

  - Foundational research in

    - Formalisms for representation of information and knowledge
    - Algorithms for reasoning with data and knowledge
    - Algorithms for knowledge acquisition

  - Applied research in

    - Semantic Web
    - Data and knowledge integration
    - Linked and Big Data
    - Ontology-based systems
    - Ontology modeling and engineering

WRIGHT STATE
UNIVERSITY

# What is this about?

- Ontology-based data integration

- Domain: geoscience, starting with ocean science

- Modular ontology engineering approach using ontology patterns.

- Aiming for flexibility and extensibility.

- As respectful as possible to individual modeling choices.

# EarthCube

- "community-driven knowledge infrastructure for geosciences"
  - well-connected environment to share data and knowledge in an open, transparent, and inclusive manner, accelerating our ability to understand and predict the Earth system
- Consists of various projects (building blocks, RCNs, SIGs) to:
  - develop key technologies,
  - promote community building,
  - explore integrative systems, and
  - prototype a governance structure.

WRIGHT STATE
UNIVERSITY

# OceanLink

- An EarthCube building block

- Applying semantic technologies for integration of existing ocean science data repositories

- Flexible, extendible, modular, respecting heterogeneity

**WRIGHT STATE**
*UNIVERSITY*

# Geosciences Data Repositories
## (a very small snapshot)

- Oceanographic data – BCO-DMO: >6000 datasets with supporting documents from 24 programs, 229 projects, 1673 deployments

- Field expeditions data – R2R: 400 expeditions per year; 3

- Conference and funded award abstracts – AGU: 30 mil. triples

- Theses, reports, journal articles – MBLWHOI Library: 5500 text documents

- Solid earth data – IEDA: hi-res bathymetry and samples from >730 cruises

- Marine geological data – IMLGS

- Ecological data – LTER

- Antarctic data – AMD

- Ocean drilling data – IODP

- Physiographic gazetteers – MRD

- ….

**October 2014 – Adila Krisnadhi**

- **Technical challenge**:
  - Lack of interoperability in terms of formats, etc.
  - Semantic heterogeneity
- **Social challenge**: Data owners/providers are <u>reluctant/unwilling</u> to participate in sharing and integration if:
  - conceptual changes have to be made to their data repositories
  - their usual business process have to be reworked, or even worse, completely discarded (note: each data repository usually represents its own research sub-community);
  - the global schema is too difficult to comprehend and manage (because the data owners are also the data consumers);
  - retrieving their own data becomes more complicated using the integrated system.

WRIGHT STATE
UNIVERSITY

- Data providers are actively involved (have a say) in the creation of the global schema.

- Definition of mappings is essentially in the hand of the data providers (knowledge engineers may help if needed, of course).

# Modular Ontology Engineering

- Model one key notion at a time

- Keep ontological commitments minimum

- Gathered constraints & requirements are formalized (e.g., with OWL) outside the modeling sessions

- Document the translation and communicate it with the domain people

- Useful if domain people can test the resulting patterns against real data

# Ontology Design Patterns

- Reusable solution to some frequently occurring ontological modeling problem emerging in different domains

- **Content pattern**: encapsulates one key notion in a particular domain, providing modular, reusable, replaceable pieces.

- By reusing generic patterns (but leaving the relationships between patterns to a specific assembly for a specific purpose), we can have a reuse while respecting heterogeneity.

- Patterns "follow" data, rather than data "follow" the patterns.

WRIGHT STATE
UNIVERSITY

# OceanLink Patterns

- Cruise

- Vessel

- Trajectory

- Person

- Organization

- Roles of Agents

- Repository Object

- Dataset

- and a few other patterns (about 15 in total)

We are not starting from zero, of course.

- Find all cruises passing through Gulf of Maine in August 2013.

- Show the trajectories of cruises in operation in September 2013.

- List all cruise vessels that departed from Woods Hole in 2012.

- Find the chief scientists of any cruise that collected samples of carbon-isotope data in Lake Superior.

- What datasets were produced by the cruise AE0901?

- Which cruises are funded by the NSF award DBI-0424599?

- List all cruises under the Ocean Flux Program.

# Cruise Pattern

# Cruise Trajectory



This reuses "Semantic Trajectory" pattern from Hu et al. COSIT 2013

DaSe Lab

$$\text{Cruise}(x) \wedge \text{providesRole}(x, y) \wedge \text{isPerformedBy}(y, z)$$
$$\wedge \text{Person}(z) \wedge \text{hasRoleType}(y, \texttt{chief\_scientist})$$
$$\rightarrow \text{hasChiefScientist}(x, z) \tag{30}$$

$$\text{Fix} \sqcap \neg \exists \text{endsAt}^-.\text{Segment} \sqsubseteq \text{StartingFix} \tag{31}$$

$$\text{Cruise}(x) \wedge \text{hasTrajectory}(x, y) \wedge \text{hasFix}(y, z) \wedge \text{StartingFix}(z)$$
$$\wedge \text{atPort}(z, p) \rightarrow \text{hasStartingPort}(x, p) \tag{32}$$

**DaSe Lab**

```
CONSTRUCT ?x rdf:type :Cruise
WHERE { ?x rdf:type r2r:Cruise. }

CONSTRUCT ?x rdf:type :Cruise
WHERE { ?x a bcodmo:Deployment;
           bcodmo:ofPlatform [a bcodmo:Vessel]. }
```

WRIGHT STATE
UNIVERSITY

- Find all ports at which the researcher "Mak Saito" stopped by in any of his expeditions.

```
DESCRIBE ?port WHERE {
    ?port a :Port.
    ?cruise :hasTrajectory ?t ;
            :hasActor ?x.
    ?t :hasFix ?f.
    ?f :atPort ?port.
    ?x rdf:type :Person; :hasLegalName "Mak Saito". }
```

WRIGHT STATE
UNIVERSITY

- Find out who joined any cruise that went through "Gulf of Maine", what their role was in the cruise, and what funding award supported their trip.

```
SELECT ?name ?role ?fund WHERE {
    ?cruise :isDescribedBy ?d; :providesRole ?r;
                :hasFix ?x.
    ?d :isFundedBy ?f.
    ?f :hasAwardID ?fund.
    ?r :hasRoleType ?role; :isPerformedBy ?p.
    ?p rdf:type :Person; :hasLegalName ?name.
    ?x :hasLocation ?pos.
    ?pl :hasSpatialFootprint ?pos; rdfs:label ?pln.
    FILTER regex(?pln, "Gulf of Maine", "i").
```

## AGU

Dr Peter Wiebe **may have** authored:

| Year | Meeting | Section | Session | Abstract |
|------|---------|---------|---------|----------|
| 2002 | Ocean Sciences | OS | OS21P | OS21P-10 |

## BCO-DMO

### Peter Wiebe was found to have the following roles

| Vessel | Cruise | Program | Role |
|--------|--------|---------|------|
| Albatross IV | AL9404 | Hydrography | Chief Scientist |
| Albatross IV | AL9508 | Hydrography | Chief Scientist |
| Albatross IV | AL9906 | Hydrography | Chief Scientist |
| Atlantis II | AT85 | Cold Core Rings | Chief Scientist |
| Endeavor | EN261 | Hydrography | Chief Scientist |
| Nathaniel B. Palmer | NBP0103 | U.S. GLOBEC Southern Ocean | Chief Scientist |
| Nathaniel B. Palmer | NBP0104 | U.S. GLOBEC Southern Ocean | Chief Scientist |
| Nathaniel B. Palmer | NBP0202 | U.S. GLOBEC Southern Ocean | Chief Scientist |
| Nathaniel B. Palmer | NBP0204 | U.S. GLOBEC Southern Ocean | Chief Scientist |
| Oceanus | OC275 | Hydrography | Chief |

- Evaluation
    - Does the pattern approach succeed in meeting both the technical and social challenges?

- Tools for assisting pattern developments
    - Ease in extending the pattern collection to cover other repositories.
    - Interesting theoretical aspect: studying various ways of ontology reuse.

- Data-to-Pattern Mappings
    - Abstraction may sometimes be more complex than the modeling on the data level, so simple query unfolding may not work.

- Reasoning
    - Entailment in queries
    - Integrity checking on data (missing or errorneous data)

# OceanLink Collaborators

- Robert Arko – Lamont-Doherty Earth Observatory, Columbia University
- Suzanne Carbotte – Lamont-Doherty Earth Observatory, Columbia University
- Cynthia Chandler – Woods Hole Oceanographic Institution
- Michelle Cheatham – Wright State University
- Timothy Finin – University of Maryland, Baltimore County
- Pascal Hitzler – Wright State University
- Krzysztof Janowicz – University of California, Santa Barbara
- Adila A. Krisnadhi – Wright State University
- Thomas Narock – Marymount University
- Lisa Raymond – Woods Hole Oceanographic Institution
- Adam Shepherd – Woods Hole Oceanographic Institution
- Peter Wiebe – Woods Hole Oceanographic Institution

# Acknowledgements

- The presented work is part of the NSF OceanLink project: "EAGER: Collaborative Research: EarthCube Building Blocks, Leveraging Semantics and Linked Data for Geoscience Data Sharing and Discovery."

# Thanks!

# References

- Aldo Gangemi. Ontology design patterns for semantic web content. ISWC 2005

- Yingjie Hu, Krzysztof Janowicz, David Carral, Simon Scheider, Werner Kuhn, Gary Berg-Cross, Pascal Hitzler, Mike Dean, and Dave Kolas. A geo-ontology design pattern for semantic trajectories. COSIT 2013.

- Willem Robert van Hage, Veronique Malaise, Roxane Segers, Laura Hollink, and Guus Schreiber. Design and use of the Simple Event Model (SEM). JWS 9(2): 2011

- Daniel Oberle, Anupriya Ankolekar, Pascal Hitzler, Philipp Cimiano, Michael Sintek, Malte Kiesel, Babak Mougouie, Stephan Baumann, Shankar Vembu, Massimo Romanelli, Paul Buitelaar, Ralf Engel, Daniel Sonntag, Norbert Reithinger, Berenike Loos, Hans-Peter Zorn, Vanessa Micelli, Robert Porzel, Christian Schmidt, Moritz Weiten, Felix Burkhardt, and Jianshen Zhou. DOLCE ergo SUMO: On Foundational and Domain Models in the SmartWeb Integrated Ontology (SWIntO). JWS 5(3): 2007

# Information representation choices

# Information representation choices

# Information Representation Choices

**R2R:**

$$Fix \sqsubseteq \exists atTime.OWL\text{-}Time{:}Temporal\ Thing \sqcap \exists hasLocation.Position$$
$$\sqcap \exists hasFix^-.SemanticTrajectory \tag{1}$$

$$Segment \sqsubseteq \exists startsFrom.Fix \sqcap \exists endsAt.Fix \tag{2}$$

$$\top \sqsubseteq\ \leq 1startsFrom.\top \tag{3}$$

$$\top \sqsubseteq\ \leq 1endsAt.\top \tag{4}$$

$$Segment \sqsubseteq \exists hasSegment^-.SemanticTrajectory \tag{5}$$

$$startsFrom^- \circ endsAt \sqsubseteq hasNext \tag{6}$$

$$hasNext \sqsubseteq hasSuccessor \tag{7}$$

$$hasSuccessor \circ hasSuccessor \sqsubseteq hasSuccessor \tag{8}$$

$$hasNext^- \sqsubseteq hasPrevious \tag{9}$$

$$hasSuccessor^- \sqsubseteq hasPredecesor \tag{10}$$

$$Fix \sqcap \neg \exists endsAt.Segment \sqsubseteq StartingFix \qquad (11)$$

$$Fix \sqcap \neg \exists startsFrom.Segment \sqsubseteq EndingFix \qquad (12)$$

$$Segment \sqcap \exists startsFrom.StartingFix \sqsubseteq StartingSegment \qquad (13)$$

$$Segment \sqcap \exists endsAt.EndingFix \sqsubseteq EndingSegment \qquad (14)$$

$$SemanticTrajectory \sqsubseteq \exists hasSegment.Segment \qquad (15)$$

$$hasSegment \circ startsFrom \sqsubseteq hasFix \qquad (16)$$

$$hasSegment \circ endsAt \sqsubseteq hasFix \qquad (17)$$

$$\exists hasSegment.Segment \sqsubseteq SemanticTrajectory \qquad (18)$$

$$\exists hasSegment^{-}.SemanticTrajectory \sqsubseteq Segment \qquad (19)$$

$$\exists hasFix.Segment \sqsubseteq SemanticTrajectory \qquad (20)$$

$$\exists hasFix^{-}.SemanticTrajectory \sqsubseteq Fix \qquad (21)$$

**DaSe Lab**

$$\text{Cruise} \sqsubseteq (=1 \ \text{hasTrajectory.Trajectory}) \tag{1}$$

$$\text{Cruise} \sqsubseteq (=1 \ \text{isUndertakenBy.Vessel}) \tag{2}$$

$$\text{Fix} \sqsubseteq \exists \text{atTime.time:TemporalEntity} \sqcap \exists \text{hasLocation.Position}$$
$$\sqcap (=1 \ \text{hasFix}^-.\text{Trajectory}) \sqcap (\leqslant 1 \ \text{nextFix.Fix}) \tag{3}$$

$$\text{Segment} \sqsubseteq (=1 \ \text{startsFrom.Fix}) \sqcap (=1 \ \text{endsAt.Fix})$$
$$\sqcap \exists \text{hasSegment}^-.\text{Trajectory} \tag{4}$$

$$\exists \text{nextFix.}\top \sqsubseteq (=1 \ \text{startsFrom}^-.\text{Segment}) \tag{5}$$

$$\exists \text{nextFix}^-.\top \sqsubseteq (=1 \ \text{endsAt}^-.\text{Segment}) \tag{6}$$

$$\text{startsFrom} \circ \text{nextFix} \sqsubseteq \text{endsAt} \tag{7}$$

**WRIGHT STATE UNIVERSITY**

$$\text{Port} \sqsubseteq \text{Place} \tag{8}$$

$$\text{Attribute}(\texttt{port\_stop\_arrival})$$

$$\text{Attribute}(\texttt{port\_stop\_departure}) \tag{9a,b}$$

$$\text{PortFix} \sqsubseteq \text{Fix} \sqcap \exists \text{atPort.Port}$$

$$\exists \text{hasAttribute}.\{\texttt{port\_stop\_arrival}\} \sqsubseteq \text{PortFix} \tag{10}$$

$$\exists \text{hasAttribute}.\{\texttt{port\_stop\_departure}\} \sqsubseteq \text{PortFix} \tag{11}$$

$$\text{atPort} \circ \text{hasSpatialFootprint} \sqsubseteq \text{hasLocation} \tag{12}$$

$$\text{hasTrajectory} \circ \text{hasSegment} \circ \text{isTraversedBy}$$
$$\sqsubseteq \text{isUndertakenBy} \tag{13}$$

$$\text{Role} \sqcap \exists \text{providesRole}^-.\text{Event} \sqsubseteq (=1 \text{ hasRoleType.RoleType})$$
$$\sqcap \exists \text{isPerformedBy.Agent} \tag{14}$$

$$\text{providesRole} \circ \text{isPerformedBy} \sqsubseteq \text{hasActor} \tag{15}$$

$$\text{Cruise} \sqsubseteq \text{Event} \tag{16}$$

$$\text{CruiseRoleType} \sqsubseteq \text{RoleType} \tag{17}$$

$$\text{CruiseRoleType}(x) \text{ for every role type } x \text{ in } (*) \tag{18a-t}$$

$$R_{\text{Cruise}} \circ \text{owl:topObjectProperty} \circ R_{\text{CruiseRoleType}}$$
$$\sqsubseteq \text{providesRoleType} \tag{19}$$

$$\text{Cruise} \equiv \exists R_{\text{Cruise}}.\text{Self}, \tag{20}$$

$$\text{CruiseRoleType} \equiv \exists R_{\text{CruiseRoleType}}.\text{Self} \tag{21}$$

**DaSe Lab**

- ## Cruise role types:
    - captain,
    - chief engineer,
    - scientist,
    - chief scientist,
    - cochief scientist,
    - postdoc scientist,
    - student,
    - graduate student,
    - undergraduate student,
    - k12 student,

- higher ed educator,
- k12 educator,
- technician,
- marine technician,
- lead marine technician,
- inspector,
- observer,
- foreign observer,
- other observer,
- scheduler,
- operator

WRIGHT STATE
UNIVERSITY

$$\text{Cruise} \sqsubseteq (=1 \text{ isDescribedBy.CruiseInformationObject}) \qquad (22)$$

$$\text{CruiseInformationObject} \sqsubseteq$$
$$(=1 \text{ hasCruiseType.CruiseType}) \qquad (23)$$

$$\text{CruiseType} \equiv \{\texttt{operational}, \texttt{transit},$$
$$\texttt{maintenance}, \texttt{other\_cruisetype}\} \qquad (24)$$

$$\text{Cruise} \sqcap \exists \text{isDescribedBy}.\exists \text{hasCruiseType}.\{\texttt{operational}\}$$
$$\equiv \exists \text{providesRole}.(\text{Role} \sqcap \exists \text{hasRoleType}.\{\texttt{chief\_scientist}\})$$
$$\sqcap \exists \text{isFundedBy.FundingAward} \qquad (25)$$

# Disjointness, Domain & Range

- Class disjointness asserted to pairs of classes, unless they are a subclass-superclass pair.

- Domain & Range use a guarded version:

$$\exists hasFix.Fix \sqsubseteq Trajectory, \quad Trajectory \sqsubseteq \forall hasFix.Fix \quad (27)$$

$$\exists hasRelatedCruiseID.rdf{:}PlainLiteral$$
$$\sqsubseteq CruiseInformationObject \quad (28)$$

$$CruiseInformationObject$$
$$\sqsubseteq \forall hasRelatedCruiseID.rdf{:}PlainLiteral \quad (29)$$

WRIGHT STATE
UNIVERSITY

# The EarthCube wish list

- Modular
- Extendible
- Sustainable
- Sliceable (you can adopt part of it without adopting all)
- Simple enough for easy adoption
- Complex enough to solve real problems
- Scalable and broad enough to cover multiple topics/domains
- Elastic/flexible enough to allow partners to decide how much they want to share
- Respectful of individual modeling choices

WRIGHT STATE
UNIVERSITY

# Example: Semantic Trajectory



Hu, Janowicz, Carral, Scheider, Kuhn, Berg-Cross, Hitzler, Dean, Kolas. COSIT 2013

# OceanLink Architecture