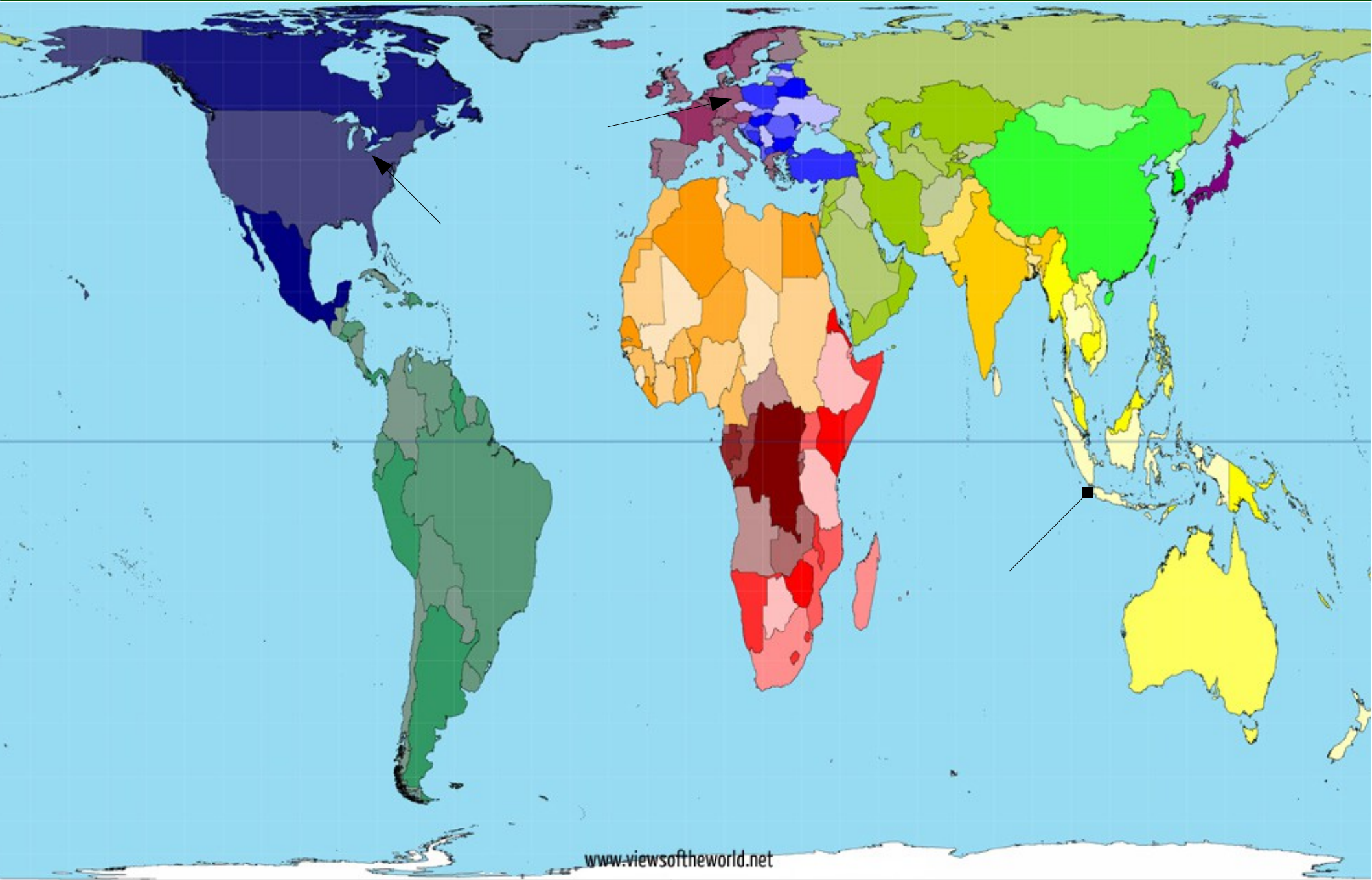


Cross-Repository Data Integration using Ontology Design Patterns

Adila Alfa Krisnadhi
DaSe Lab for Data Semantics
Wright State University



- 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

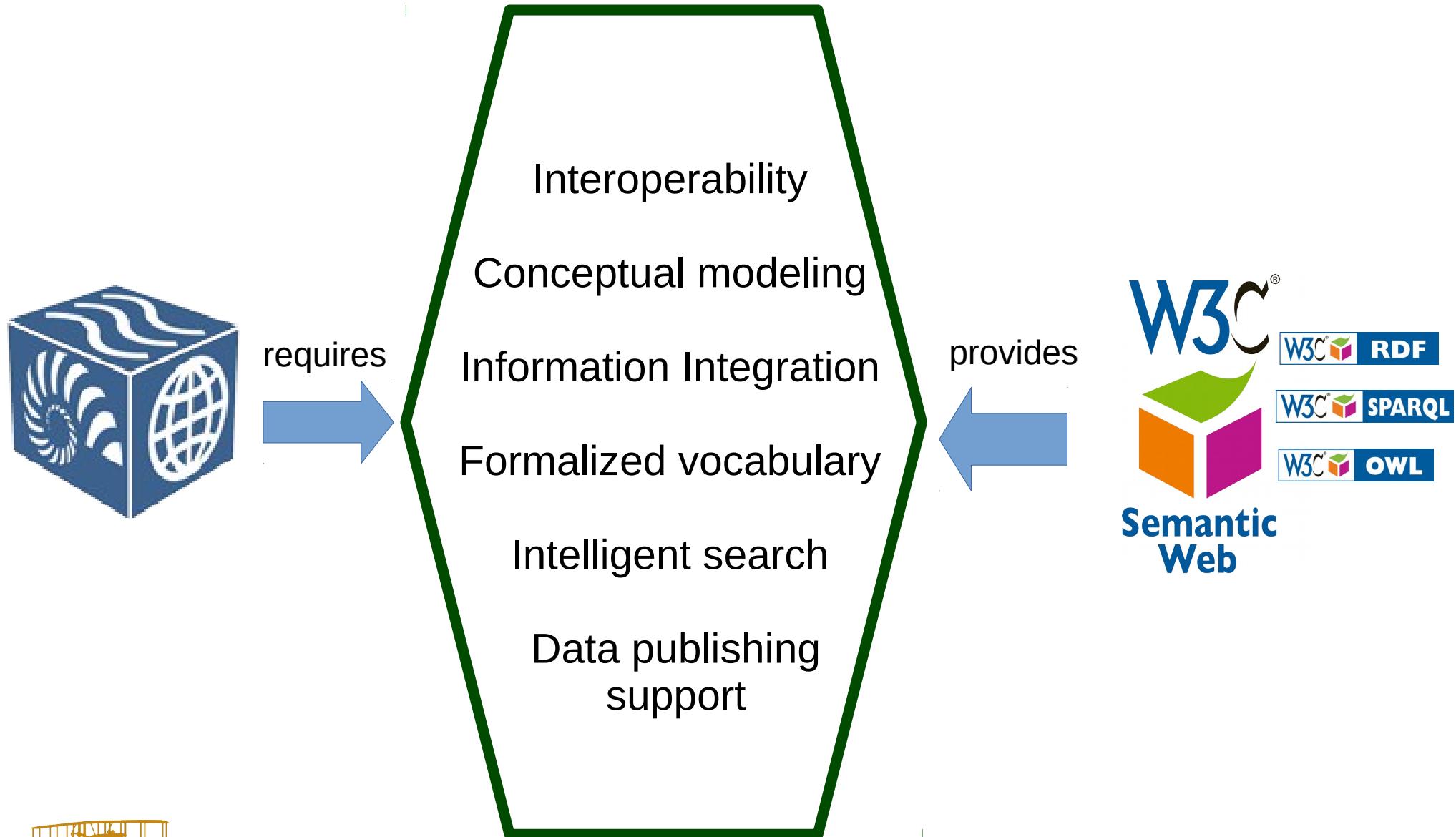
What is this about?

- 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.

Background

- “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.





- An EarthCube building block
- Applying semantic technologies for integration of existing ocean science data repositories
- Flexible, extendible, modular, respecting heterogeneity

NSF award 1354778 "EAGER: Collaborative Research: EarthCube Building Blocks, Leveraging Semantics and Linked Data for Geoscience Data Sharing and Discovery."

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
-

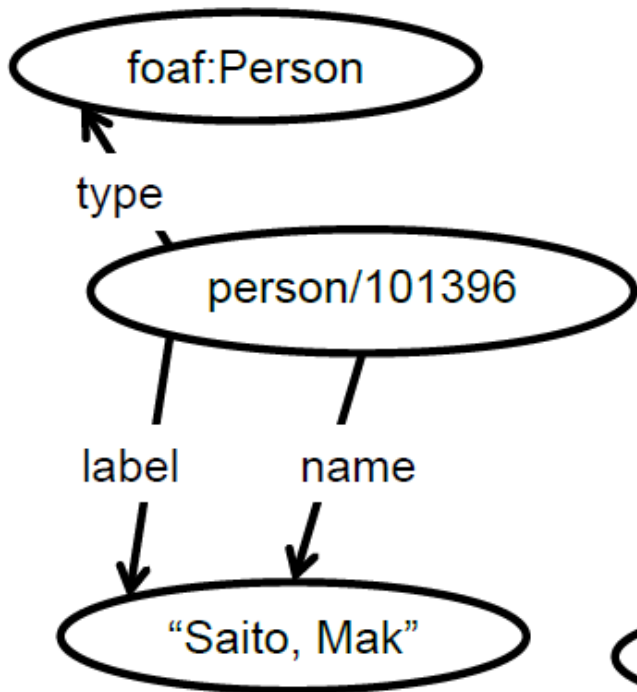


- **Technical challenge:**
 - Lack of interoperability in terms of formats, etc.
 - Semantic and content heterogeneity
- **Social challenge:** Data owners/providers are reluctant/unwilling to participate in sharing and integration if:
 - conceptual changes have to be made to their data repositories
 - 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.

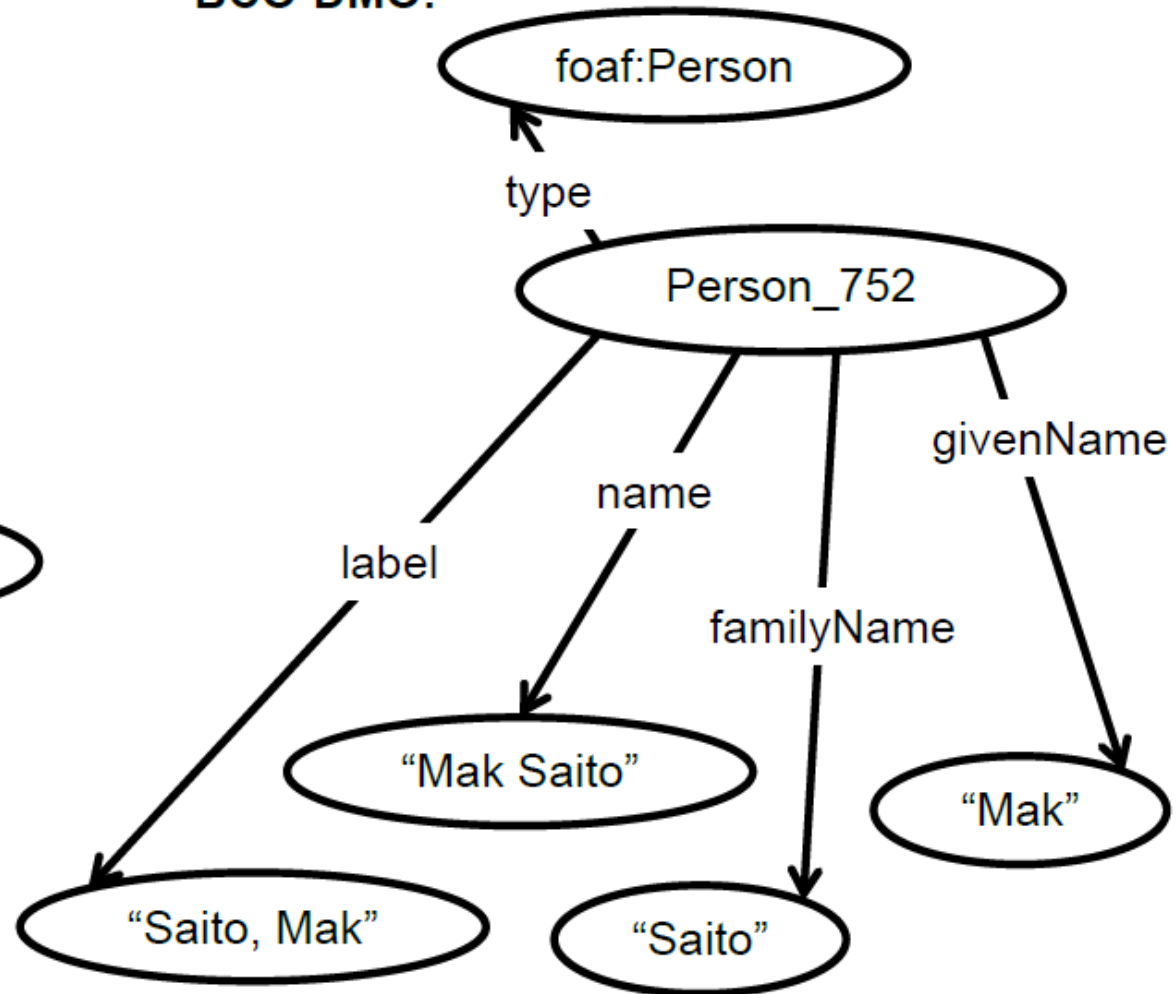


Different ways of modeling Person

R2R:

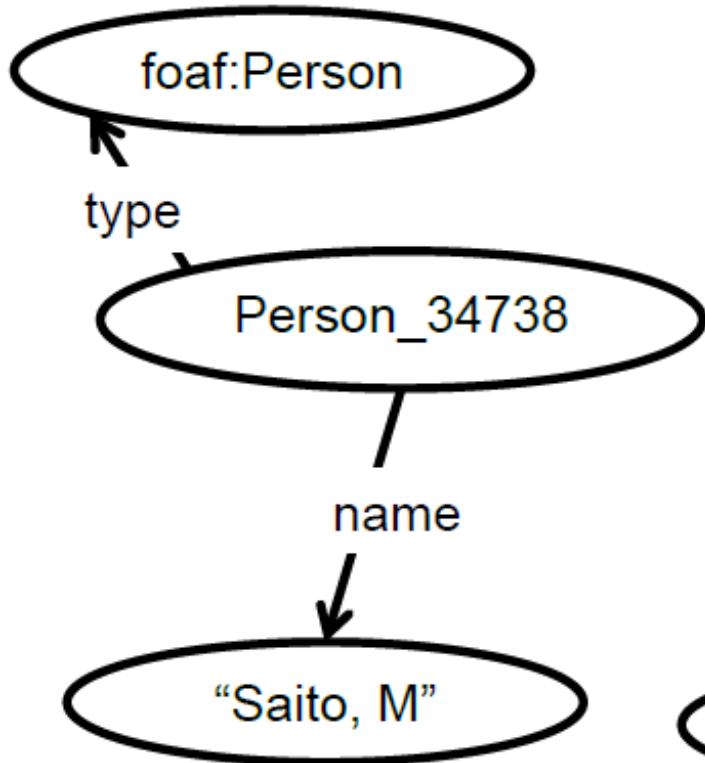


BCO-DMO:

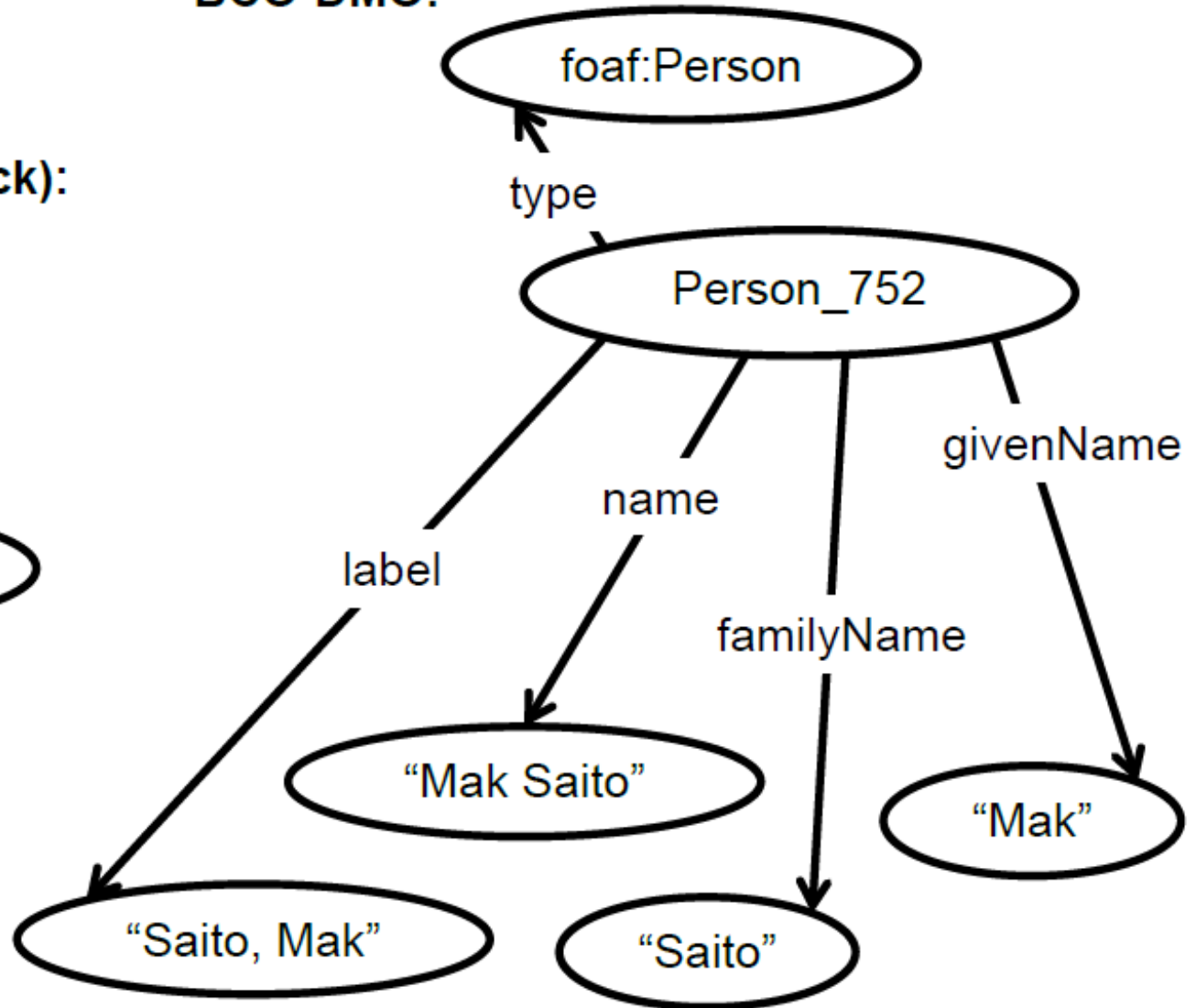


Different ways of modeling Person

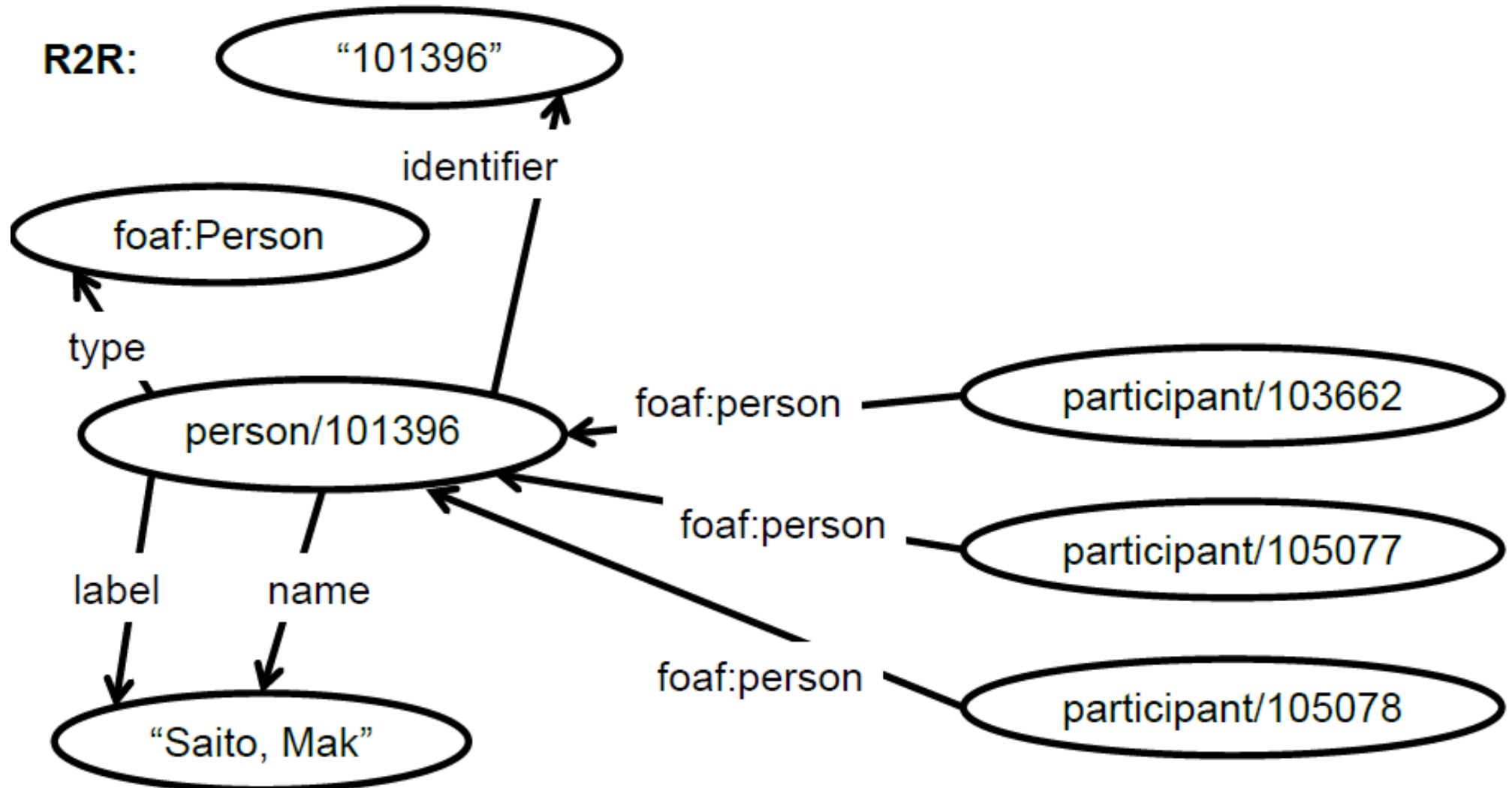
AGU Abstracts (Tom Narock):



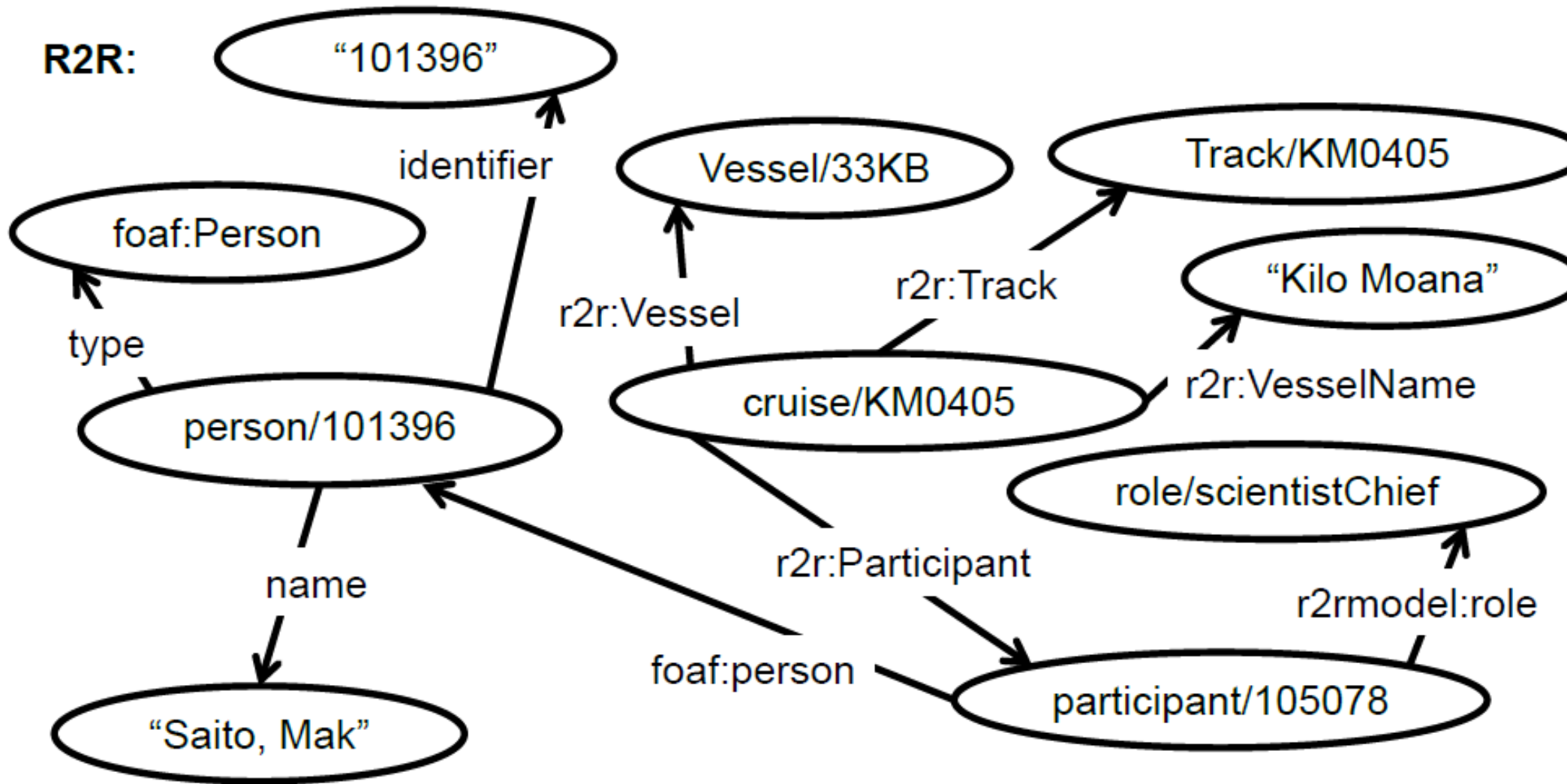
BCO-DMO:



Different ways of modeling Person



Different ways of modeling Person



Ontology (Design Pattern) Based Data Integration

- Data Integration Problem:
 - The problem of providing a user with a unified view over data residing at different, autonomous, and possibly heterogeneous data sources.
 - The unified view is called the *global schema*.
 - When a user issues a query over the global schema, the system should translate it into queries over suitable data sources and assemble the results into the final answer to be presented to the user.

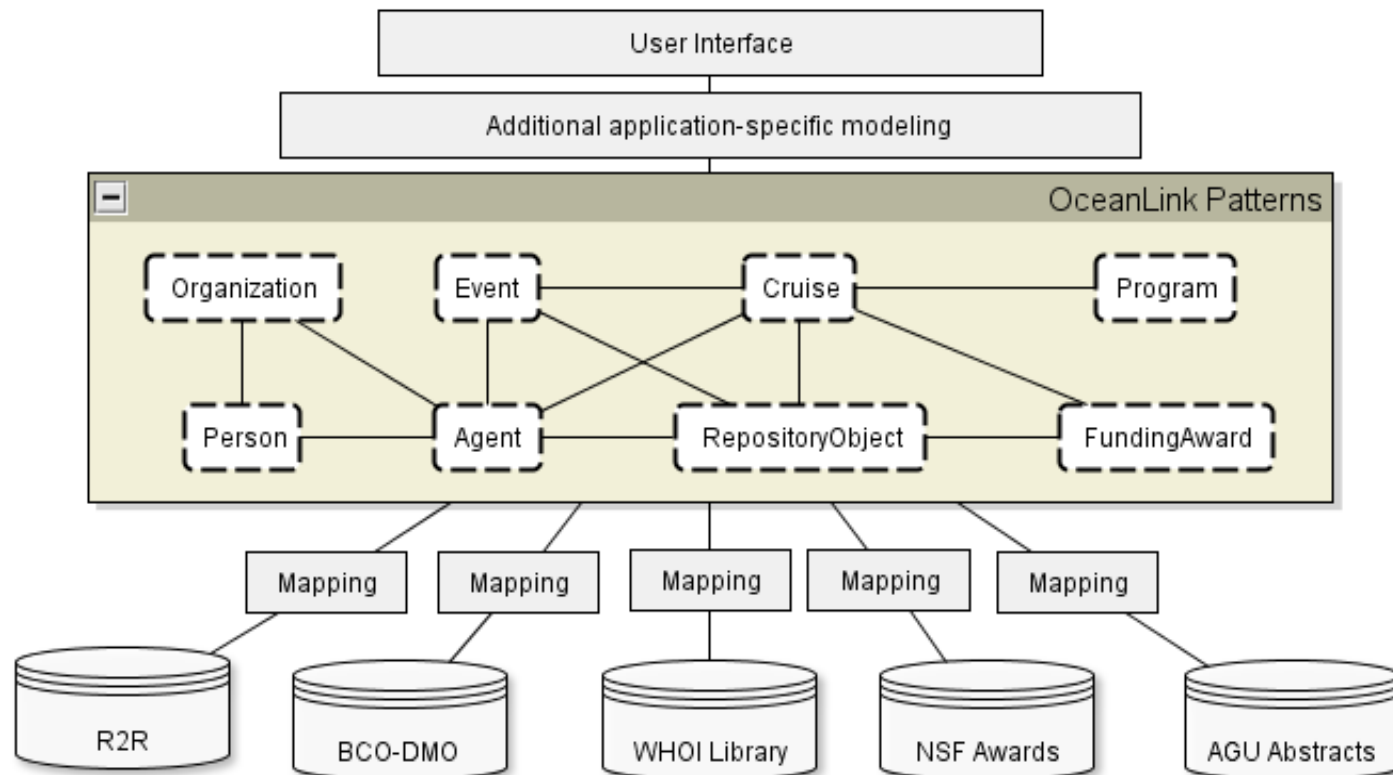
- Interoperability through RDF standard:
 - Web standards on data model (RDF: a set of triples), querying language and protocol (SPARQL; REST)
 - Each data repository only needs to make their data available for SPARQL querying.
- Vocabulary of global schema is defined using ontology.
 - Allows more expressive conceptualization using a language like OWL – ontology given as a set of logical axioms.
 - Bridging semantic heterogeneity
 - *But, using overarching monolithic ontology is too difficult, complex, unwieldy, not flexible, too constraining, not easily extendible, etc.*



- Model one key notion at a time
- Keep ontological commitments minimum (avoid too constraining axiomatic statements)
- 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 experts can test the resulting patterns against real data

- 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.

- 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).



- BCO-DMO
 - Award, Cruise, Device, Fileset, Format, Geometry, Holding, Model, Organization, Participant, Person, Port, Product, Program, Report, Repository, Track, Vessel
- R2R
 - Program, Project, Deployment, Dataset, Instrument, Parameter, People, Affiliation, Funding Source, Award
- AGU Abstract
 - Meeting, Meeting Section, Meeting Session, Abstracts of Contribution (Poster, Talk, etc.), People, Organization
- MBLWHOI Library
 - Article, Dataset, Book, Report, Person, Organization
- NSF Funded Award Repository
 - Award, Person, Organization, Publication, Dataset

- 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.

Engineering an Ontology Design Pattern

- Find all cruises passing through Gulf of Maine in August 2013.
- List all cruise vessels that departed from Woods Hole in 2012.
 - Cruise is conducted on some vessel, running on some track (trajectory), which is bounded by some spatio-temporal boundary.
- Find the chief scientists of any cruise that collected samples of carbon-isotope data in Lake Superior.
 - Some people/organization performs certain roles in a cruise.
- What datasets were produced by the cruise AE0901?
 - Cruise has some (unique?) identifier, name, etc
 - Cruise (activities) may result in some reports, datasets, etc.
- Which cruises are funded by the NSF award DBI-0424599?
- List all cruises under the Ocean Flux Program.
 - Cruise is funded by some funding award or part of some program.

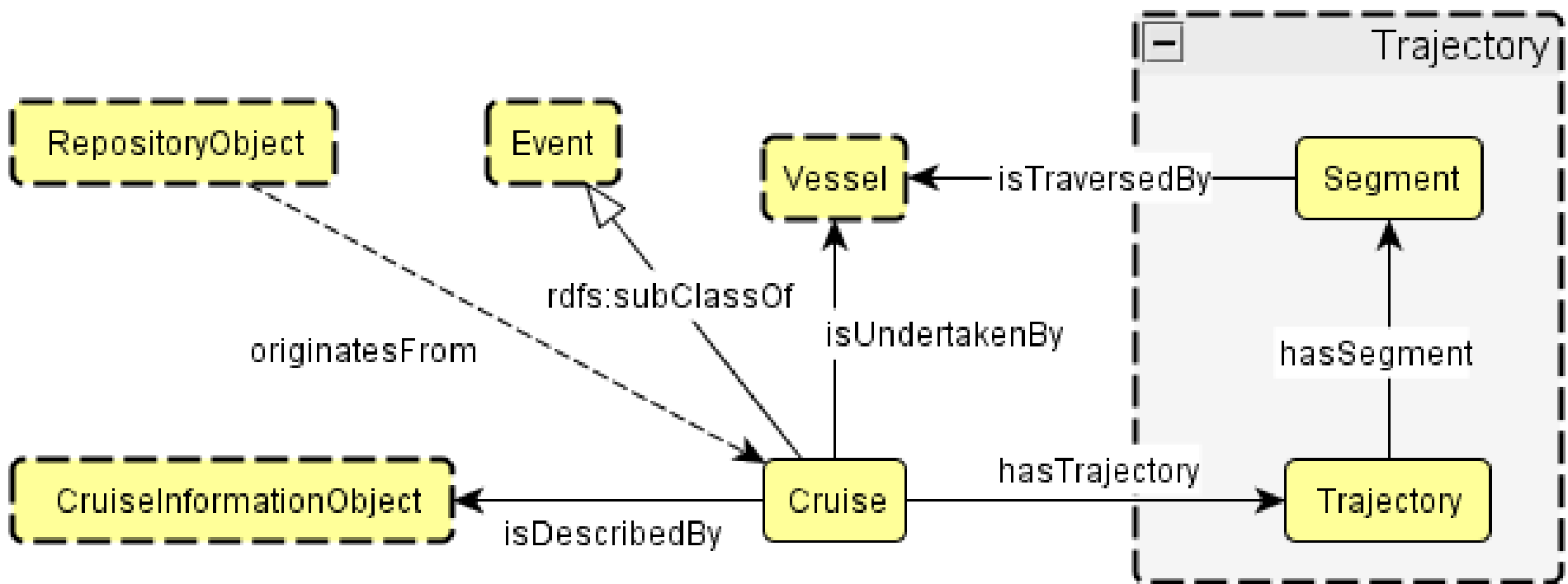


Can we reuse existing patterns?

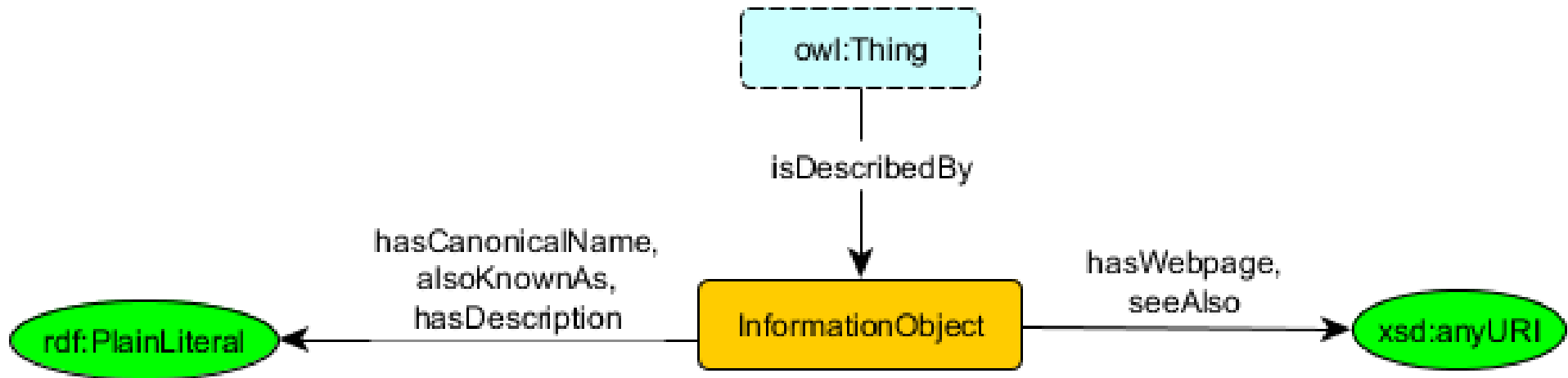
- Some of the use cases given by the CQs fit some existing patterns.
 - Simple Event Model [van Hage, et al., JWS 2011]; **lack of formalization**
 - Semantic Trajectory [Hu, et al., COSIT 2013]; **formalization in OWL, but some details need to be modified**
 - Information Object from DOLCE [Oberle, et al., JWS 2007]; **pattern only used as an inspiration, formalization is redone for our purpose**
- Reuse is not simply done via ontology import, but requires adjustment and massaging to make it appropriate for our needs.
- The challenge is finding appropriate existing patterns for reuse, and performing the suitable adjustments



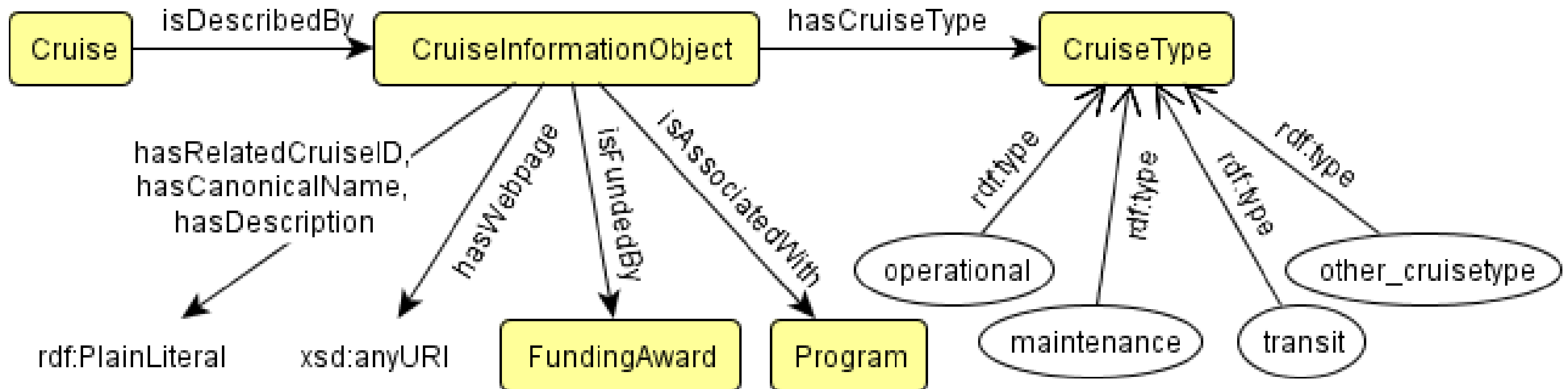
Cruise Pattern



InformationObject pattern



Cruise Information Object



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

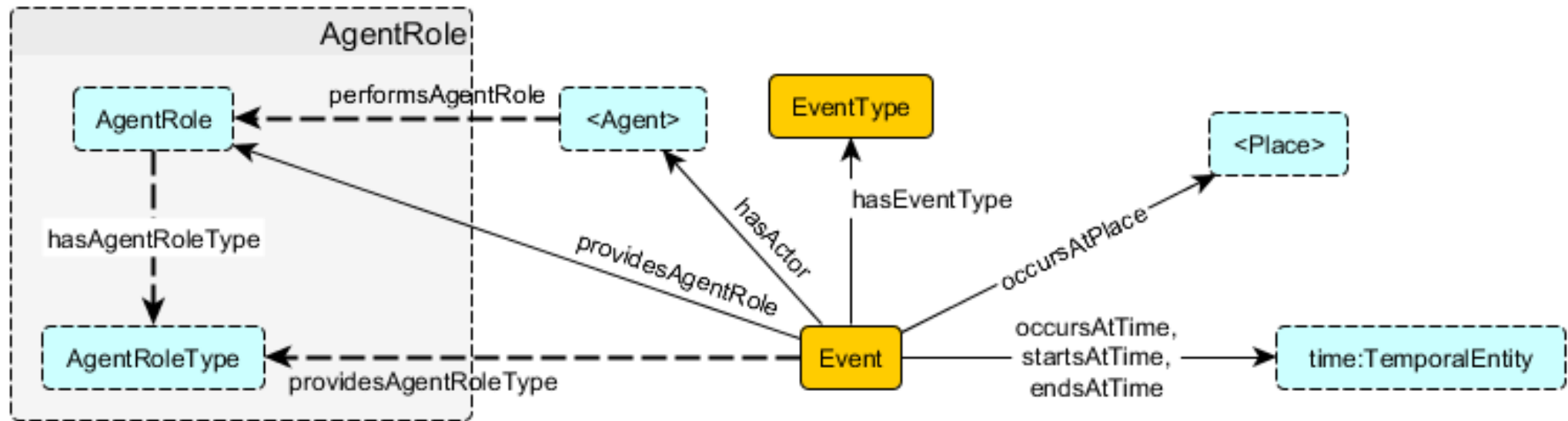
$$\begin{aligned} \text{CruiseInformationObject} \sqsubseteq \\ (=1 \text{ hasCruiseType.CruiseType}) \end{aligned} \quad (23)$$

$$\begin{aligned} \text{CruiseType} \equiv \{ \text{operational, transit,} \\ \text{maintenance, other_cruisetype} \} \end{aligned} \quad (24)$$

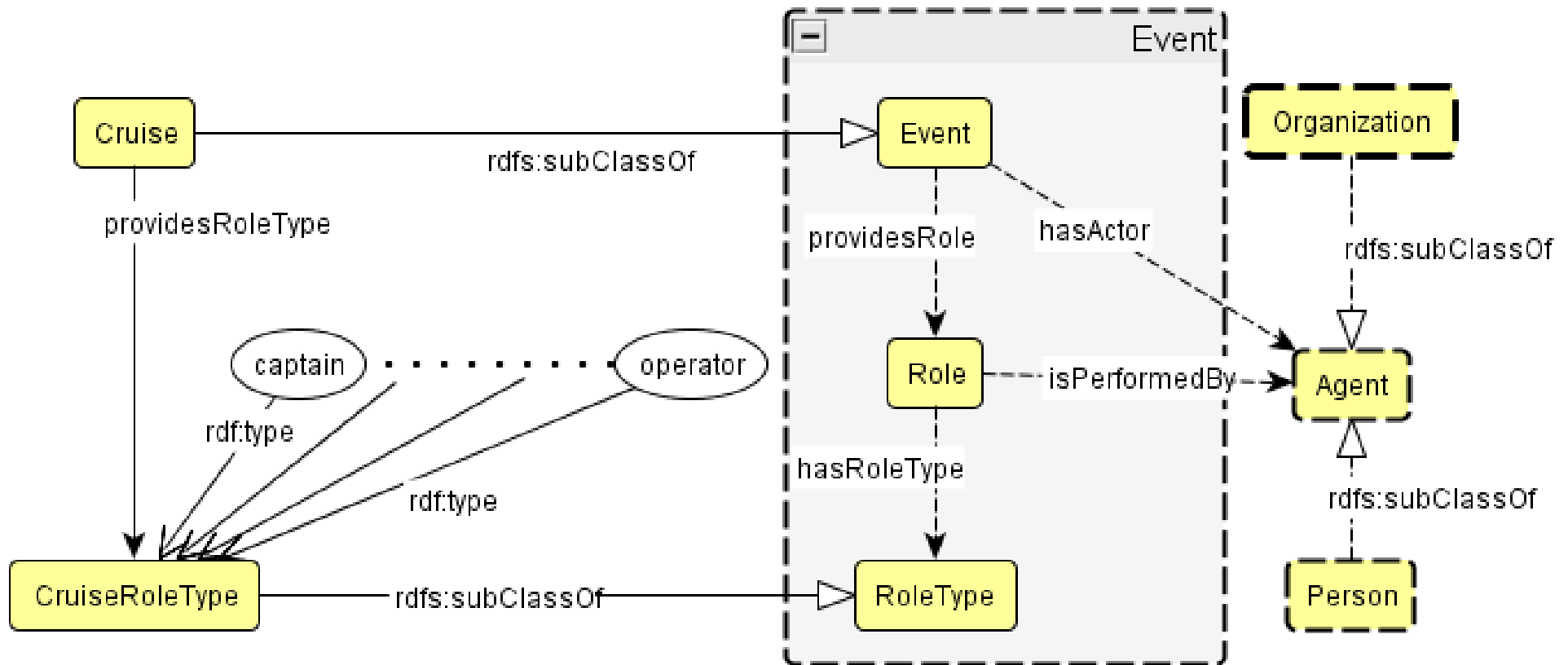
$$\begin{aligned} \text{Cruise} \sqcap \exists \text{isDescribedBy.} \exists \text{hasCruiseType.} \{ \text{operational} \} \\ \equiv \exists \text{providesRole.} (\text{Role} \sqcap \exists \text{hasRoleType.} \{ \text{chief_scientist} \}) \\ \sqcap \exists \text{isFundedBy.FundingAward} \end{aligned} \quad (25)$$



Event pattern



Roles (Cruise as Event)



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

$$\text{providesRole} \circ \text{isPerformedBy} \sqsubseteq \text{hasActor} \quad (15)$$

$$\text{Cruise} \sqsubseteq \text{Event} \quad (16)$$

$$\text{CruiseRoleType} \sqsubseteq \text{RoleType} \quad (17)$$

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

$$R_{\text{Cruise}} \circ \text{owl:topObjectProperty} \circ R_{\text{CruiseRoleType}} \sqsubseteq \text{providesRoleType} \quad (19)$$

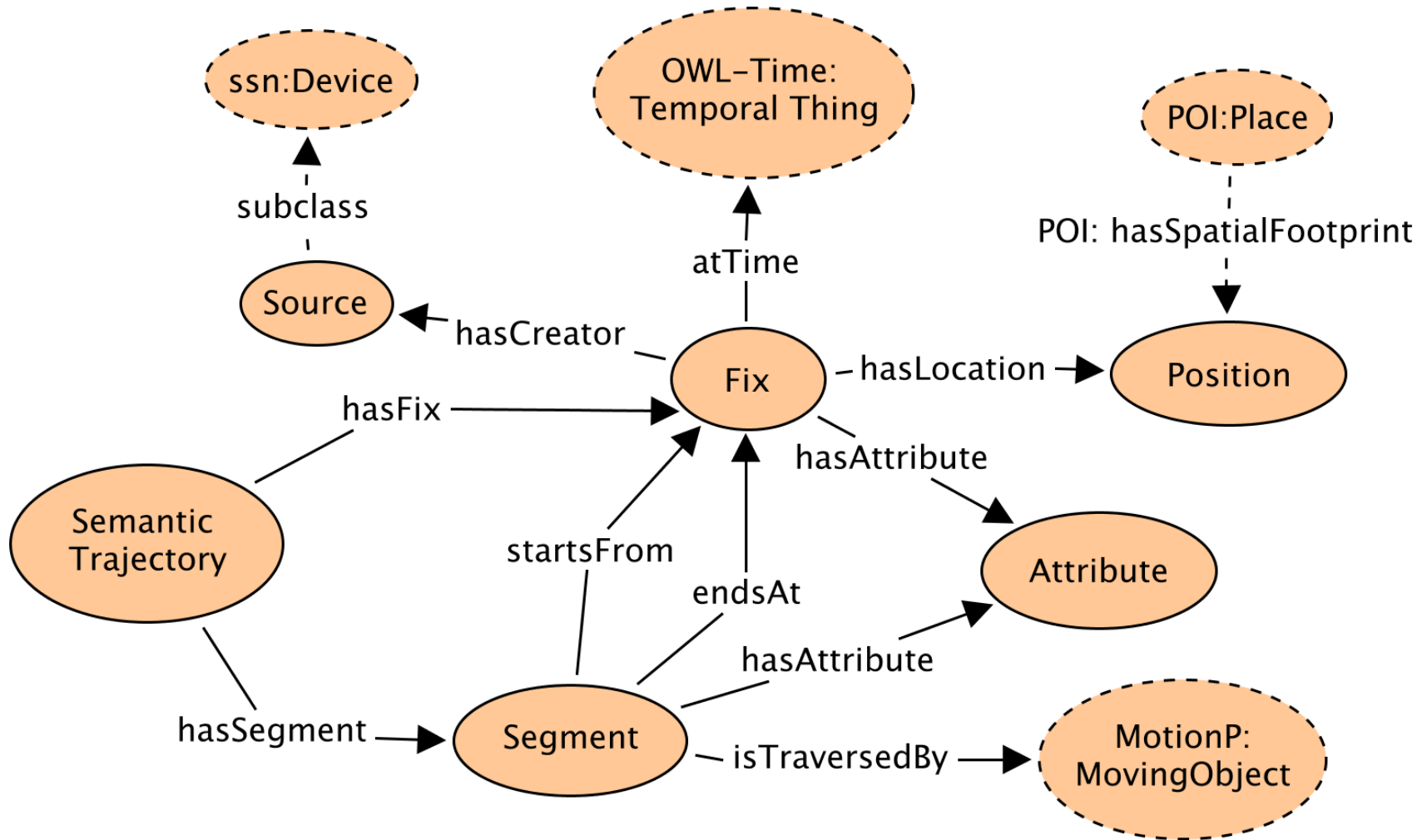
$$\text{Cruise} \equiv \exists R_{\text{Cruise}} . \text{Self}, \quad (20)$$

$$\text{CruiseRoleType} \equiv \exists R_{\text{CruiseRoleType}} . \text{Self} \quad (21)$$



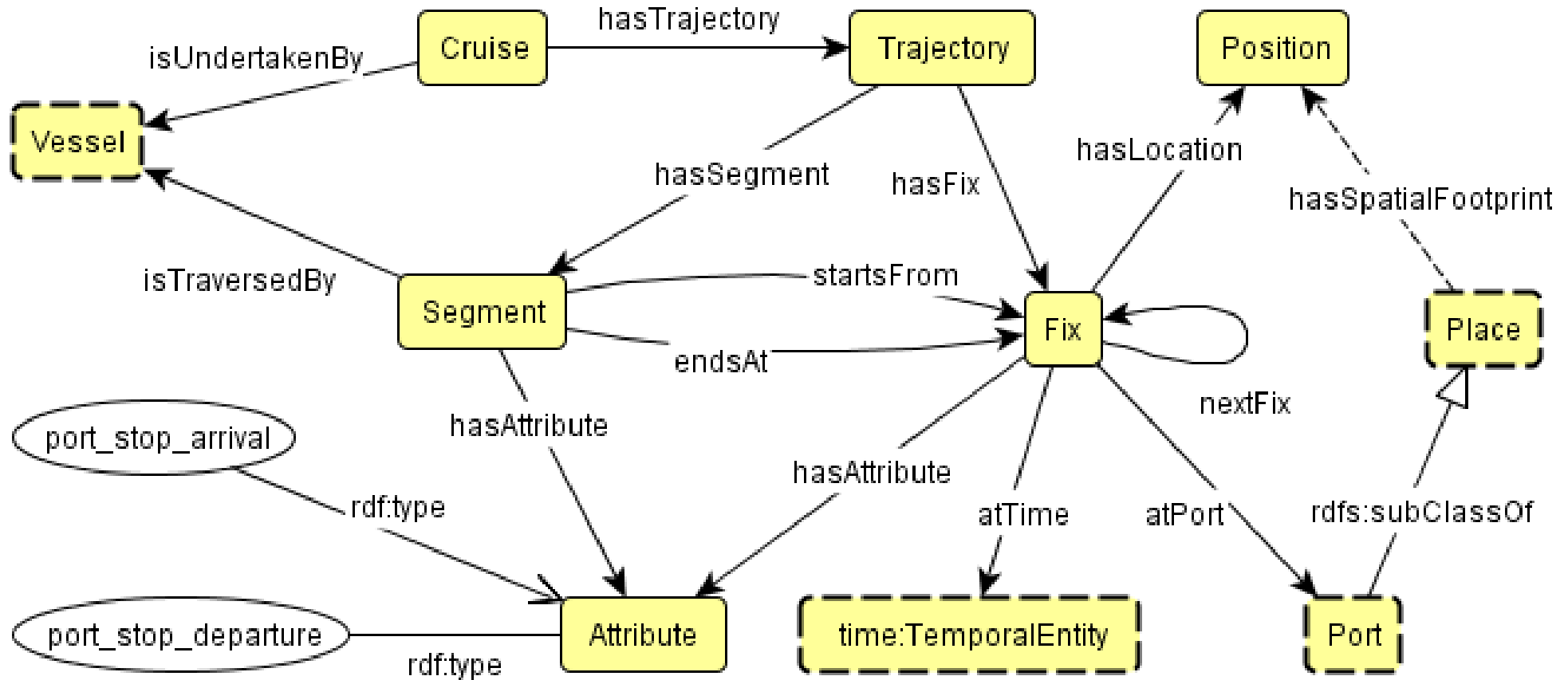
- Cruise role types:
 - captain,
 - chief engineer,
 - scientist,
 - chief scientist,
 - cochief scientist,
 - postdoc scientist,
 - student,
 - graduate student,
 - undergraduate student,
 - k12 student,
 - higher ed educator,
 - k12 educator,
 - technician,
 - marine technician,
 - lead marine technician,
 - inspector,
 - observer,
 - foreign observer,
 - other observer,
 - scheduler,
 - operator

Semantic Trajectory Pattern



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

Cruise Trajectory



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

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

$$\begin{aligned} \text{Fix} \sqsubseteq \exists \text{atTime.time:TemporalEntity} \sqcap \exists \text{hasLocation.Position} \\ \sqcap (=1 \text{ hasFix}^- \text{.Trajectory}) \sqcap (\leq 1 \text{ nextFix.Fix}) \end{aligned} \quad (3)$$

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

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

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

$$\text{startsFrom} \circ \text{nextFix} \sqsubseteq \text{endsAt} \quad (7)$$



Port \sqsubseteq Place (8)

Attribute(port_stop_arrival)

Attribute(port_stop_departure) (9a,b)

PortFix \sqsubseteq Fix \sqcap \exists atPort.Port

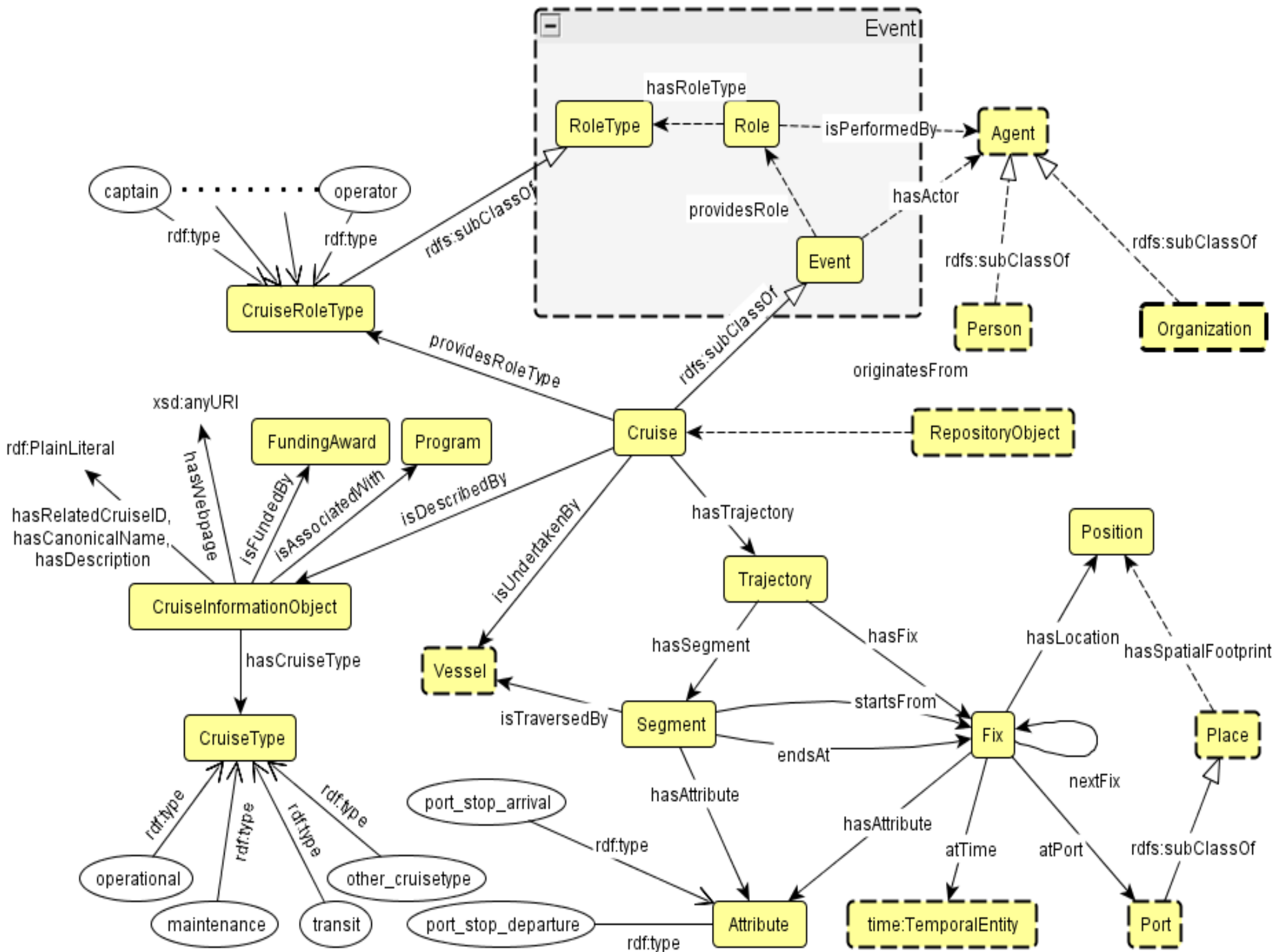
\exists hasAttribute.{port_stop_arrival} \sqsubseteq PortFix (10)

\exists hasAttribute.{port_stop_departure} \sqsubseteq PortFix (11)

atPort \circ hasSpatialFootprint \sqsubseteq hasLocation (12)

hasTrajectory \circ hasSegment \circ isTraversedBy
 \sqsubseteq isUndertakenBy (13)





- Class disjointness asserted to pairs of classes, unless they are a subclass-superclass pair.
- Domain & Range use a guarded version:

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

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

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

Queries, Query Shortcuts, and Mappings

- 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". }
```

- 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").
```



$$\begin{aligned} & \text{Cruise}(x) \wedge \text{providesRole}(x, y) \wedge \text{isPerformedBy}(y, z) \\ & \wedge \text{Person}(z) \wedge \text{hasRoleType}(y, \text{chief_scientist}) \\ & \quad \rightarrow \text{hasChiefScientist}(x, z) \end{aligned} \quad (30)$$

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

$$\begin{aligned} & \text{Cruise}(x) \wedge \text{hasTrajectory}(x, y) \wedge \text{hasFix}(y, z) \wedge \text{StartingFix}(z) \\ & \quad \wedge \text{atPort}(z, p) \rightarrow \text{hasStartingPort}(x, p) \end{aligned} \quad (32)$$



Mapping Rules (Examples)

For R2R

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

For BCO-DMO

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

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

- Full-fledged implementation (current prototype: www.oceanlink.org)
- Evaluation (w.r.t. 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.
- Mappings
 - Abstraction may sometimes be more complex than the modeling on the data level, so simple query unfolding may not work (need to introduce blank nodes)
- Reasoning
 - Entailment in queries
 - Co-reference resolution
 - Integrity checking on data (missing or erroneous data)

- 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

- 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.”

- 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



Thanks!