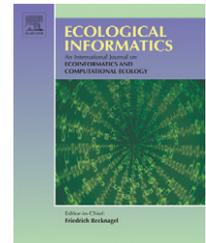


available at www.sciencedirect.comwww.elsevier.com/locate/ecolinf

An ontology for describing and synthesizing ecological observation data

Joshua Madin^{a,e,*}, Shawn Bowers^b, Mark Schildhauer^a, Serguei Krivov^c,
Deana Pennington^d, Ferdinando Villa^c

^aNational Center for Ecological Analysis and Synthesis, University of California, Santa Barbara, California 93101, USA

^bGenome Center, University of California, Davis, California 95616, USA

^cGund Institute for Ecological Economics, University of Vermont, Burlington, Vermont 05405, USA

^dUniversity of New Mexico, Albuquerque, New Mexico 87131, USA

^eDepartment of Biological Sciences, Macquarie University, New South Wales 2109, Australia

ARTICLE INFO

Article history:

Received 6 February 2007

Accepted 30 May 2007

Keywords:

Ecology

Observation

Measurement

Ontology

Data discovery

Data integration

ABSTRACT

Research in ecology increasingly relies on the integration of small, focused studies, to produce larger datasets that allow for more powerful, synthetic analyses. The results of these synthetic analyses are critical in guiding decisions about how to sustainably manage our natural environment, so it is important for researchers to effectively discover relevant data, and appropriately integrate these within their analyses. However, ecological data encompasses an extremely broad range of data types, structures, and semantic concepts. Moreover, ecological data is widely distributed, with few well-established repositories or standard protocols for their archiving and retrieval. These factors make the discovery and integration of ecological data sets a highly labor-intensive task. Metadata standards such as the Ecological Metadata Language and Darwin Core are important steps for improving our ability to discover and access ecological data, but are limited to describing only a few, relatively specific aspects of data content (*e.g.*, data owner and contact information, variable “names”, keyword descriptions, etc.). A more flexible and powerful way to capture the semantic subtleties of complex ecological data, its structure and contents, and the inter-relationships among data variables is needed.

We present a formal ontology for capturing the semantics of generic scientific observation and measurement. The ontology provides a convenient basis for adding detailed semantic annotations to scientific data, which crystallize the inherent “meaning” of observational data. The ontology can be used to characterize the context of an observation (*e.g.*, space and time), and clarify inter-observational relationships such as dependency hierarchies (*e.g.*, nested experimental observations) and meaningful dimensions within the data (*e.g.*, axes for cross-classified categorical summarization). It also enables the robust description of measurement units (*e.g.*, grams of carbon per liter of seawater), and can facilitate automatic unit conversions (*e.g.*, pounds to kilograms). The ontology can be easily extended with specialized domain vocabularies, making it both broadly applicable and highly customizable. Finally, we describe the utility of the ontology for enriching the capabilities of data discovery and integration processes.

Published by Elsevier B.V.

* Corresponding author.

E-mail address: madin@nceas.ucsb.edu (J. Madin).

1. Introduction

Ecology is an inherently multidisciplinary science that explores how physical and biological factors and their inter-relationships establish the structure and function of living systems. Accordingly, the range of data that can inform ecological analyses is incredibly broad, often involving perspectives from many fields in the earth sciences (*e.g.*, geography, oceanography and hydrology) and life sciences (*e.g.*, genetics and physiology). The need to access these diverse data sources becomes especially acute when undertaking *synthetic* analyses to address broad ecological questions, such as the impacts of deforestation on the global balance of greenhouse gases, or the link between biodiversity losses and the productivity of our world's fisheries. Ecological insights are critical to understanding many complex real world issues that have vital implications for the quality and sustainability of life on this planet.

Approaches to ecological synthesis today leverage the rapidly growing amounts of data available through the Internet. However, current methods for finding and interpreting potentially relevant data are extremely primitive and inefficient, which severely impedes progress in accomplishing *synthetic* ecological science (Pickett *et al.*, 1994). The lack of advanced technical tools for data exploration and interpretation has long been recognized (Chen, 1976; Batini *et al.*, 1992; Hammer and Macleod, 1999), and proposed enhancements are still largely lacking in practical, non-proprietary implementations.

Effective data discovery is particularly problematic in ecology, where traditionally small, focused studies employed largely *ad hoc* data management solutions, often consisting of flat files or spreadsheets with minimal formal structure and little to no metadata documentation. This situation was viable when researchers worked only with their own data, and data management was considered merely a provisional framework for accomplishing some specific analyses, after which, one moved on to other research questions and data analyses. Researchers maintained many of the details of their data in their memory, with maximum cognizance of the relevant subtleties and issues in the data ideally occurring simultaneously with the period when they were actively being analyzed (Michener *et al.*, 1997).

Recently there has been a growing recognition of the need to both preserve ecological data after their intended usage was completed, as well as to extend data collection events through time to discern long-term trends in ecological processes through intensive site-based studies (Michener, 2000). This recognition raised concerns about the lack of protocols and services for preserving ecological data (Gross and Pake, 1995), and clarified the need for reducing the possibilities of “data entropy” (Michener *et al.*, 1997). Efforts grew to develop metadata standards that can systematically structure the types of information that researchers should document about their data, making these far more effective for informing future studies (Jones *et al.*, 2006).

The *Ecological Metadata Language (EML)* was developed through a community-based effort involving researchers and information managers from several institutions charged with accomplishing ecological synthesis and long-term research (Jones *et al.*, 2001). While intended primarily for the purpose of

preserving critical metadata about ecological data sets, it is essentially a generic standard for describing tabular data, in addition to a number of other data formats (Fegraus *et al.*, 2005). While EML is a growing standard for data documentation in the ecological field, practical experience using this standard has revealed that metadata alone has some serious shortcomings in terms of the capabilities it can provide scientists in data discovery and interpretation. These shortcomings are particularly severe in ecology due to the heterogeneity of topics studied, and the relative lack of standardized protocols and methods when measuring variables of ecological interest.

Metadata languages have also been developed for describing natural history specimen data, such as Darwin Core (Darwin Core, 2004). Both the Darwin Core and EML metadata standards primarily focus on describing data format (*i.e.*, describing data structure) along with high-level contextual information (often by adopting Dublin-core style attributes such as who created a data set and when [DCMI, 2006]). Furthermore, these standards generally lack support for capturing even the basic “semantics” of data—*i.e.*, information that broadens the capability for understanding or interpreting the content and relevance of the data from a disciplinary perspective. For example, while EML allows one to declare that a data set contains an attribute labeled “biom” (*e.g.*, referring to a biomass measure), it is not possible using EML to determine if it is compatible with another attribute labeled “weight” or “kg”. What is needed is a way to capture such concepts and relationships in formal models that can then be used to draw logical conclusions (*e.g.*, consistency, equivalence) without human intervention. The creation of such a framework must also address the need for a simple mechanism to assist scientists in mapping their observations onto such models.

This paper describes an approach for enhancing the capability of ecological scientists to more powerfully discover, interpret, and reuse data in support of *synthetic* research. While provision of access to “others” data raises some interesting issues and challenges with regards to the sociology of data-sharing and intellectual property rights, the focus here is solely on addressing several of the most pressing technological impediments to accomplishing scientific synthesis: *data discovery* and (legitimate) *integration*. Discovery is the process of locating relevant and available data related to a specified topic of interest. This process is currently hampered by the lack of well-described data to begin with, and compounded by the inability to clearly explicate and explore basic semantic notions within and across data sets (*e.g.*, that biomass is a weight and that “dry weight” is a biomass and a weight). Integration is the process of merging compatible data once these are discovered. Here, we present a formal *ontological* framework for capturing the essential semantic information of observational data sets to better facilitate the discovery and integration of ecological information, thus aiding ecologists in synthesizing knowledge for answering larger ecological questions.

Ontologies are representations of the knowledge within a domain of interest, defined *via* the terminology (concepts) used within the domain and the properties and relationships among domain objects (Baader *et al.*, 2003). In this way, ontologies represent one enabling mechanism for providing

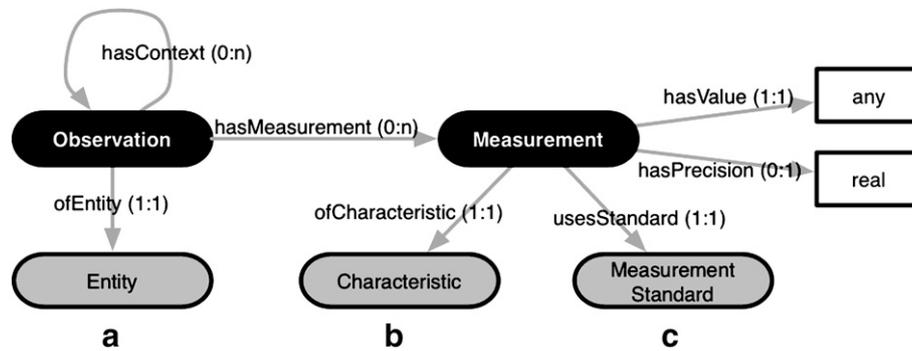


Fig. 1 – The core classes (ellipses) and properties (arrows) of the Extensible Observation Ontology (OBOE). Each Observation is of some Entity, and can provide context for the Observation of another Entity. A Characteristic of an Entity can be represented through a Measurement. Measurements relate Characteristics to a Measurement Standard via a Value and, if applicable, a Precision. Measurements are taken by a Recorder (human or non-human) using a Protocol at a particular Time and Place (shaded properties, see text for details). Observations may have multiple Measurements. Entity (a), Characteristic (b), and Measurement Standard (c) (shaded classes) provide extension points for domain-specific ontologies (see Fig. 3). Numbers in parentheses denoted min:max cardinality for properties.

more comprehensive data discovery and integration (Jones et al., 2006). As a simple example, if instances of *biomass* are defined as instances of *weight* in a particular domain ontology, then data about *biomass* will be discovered when searching for data about *weight* (i.e., based on an ontology reasoning system [Baader et al., 2003]). Moreover, these data are compatible through being semantically classified as *weights*, and can potentially be merged. A number of formal (logic-based) languages exist to capture ontologies, including description logics (Baader et al., 2003), semantic networks (Sowa, 1999), and the more recent RDF and OWL web-based standards (McGuinness and van Harmelen, 2004). While ontologies are being used successfully in a number of biological and medical informatics projects (The Gene Ontology Consortium, 2000; Rosse and Mejino, 2003; Bard and Rhee, 2004), widespread support for ontology-based approaches and ontology development has yet to be adopted within the field of ecology.

Our work on ontologies is within the context of the Science Environment for Ecological Knowledge (SEEK; <http://seek.ecoinformatics.org>) project, which aims to develop technology to discover, access, integrate, and analyze distributed ecological information (e.g., using scientific workflows [Ludäscher et al., 2006; Kepler Project, 2006]). The project's approach is to extend EML to support the semantic annotation of ecological data sets, such that EML data-set descriptions can use terms drawn from OWL-DL ontologies. A benefit of using OWL-DL (which itself is based on description logic) is that it supports a “natural” representation for formalizing terms. In particular, named classes (such as *Biomass* or *Plot*) define sets in OWL-DL, where each member of a set is considered an “instance” of the corresponding class (e.g., the class *Biomass* might denote the set of all biomasses; *Plot* the set of all physical plots). Class definitions are typically intentional, i.e., classes are defined based on their name or possibly other constraints, without enumerating their instances. The “intent” of a class is usually further elaborated by relating it to other classes. The “is-a” relationship defines class specializations as subsets. For example, the expression “*Biomass* is-a *Weight*” implies all biomass instances are also valid weight instances. OWL-DL

also supports user-defined properties, e.g., to capture part-of relationships between classes, as well as cardinality and other (set-based) constraints (Baader et al., 2003).

While formal languages such as OWL-DL provide a means to capture ontologies, the quality of the realized ontology will determine its utility for assisting in data discovery and integration. Additionally, as the number of ontologies and their included terms increases, organizing these into a coherent framework becomes increasingly complex, as recognized within the biological community, e.g., see [Bard and Rhee, 2004]. In this paper, we describe the SEEK Extensible Observation Ontology (OBOE), which aims at providing a core ontology framework for semantically annotating observational data sets. Our framework defines a formal ontology based on the concepts of *Observation*, *Measurement*, (Ecological) *Entity*, *Characteristic*, and *Measurement Standard* (e.g., physical units) (Fig. 1), providing a structured yet generic approach for semantic data annotation and for developing (and combining) domain-specific ecological ontologies. Our approach differs from other ontology-based descriptions of ecological environmental information (e.g., Keet, 2005; Smith and Varzi, 1999a,b; Smith, 2001; Gruber and Olsen, 1994; Brillhante, 2003; Cox, 2006; Williams et al., 2006; Schnetz and Mirtl, 2003) in that we focus specifically on providing: (i) a robust framework for describing generic scientific observations; (ii) a structured approach for easily building and sharing domain-specific ontology extensions; and (iii) data discovery and integration services, via semantic annotations to the ontology, across varied ecological observation data (and not just for a specific, specialized domain). In Section 2, we describe our framework using a number of real-world examples, and illustrate how it can be extended with domain-specific ontologies. The ontology framework presented here has evolved through various earlier efforts within SEEK to develop formal ecological ontologies (Berkley et al., 2005; Bowers et al., 2005; Bowers and Ludäscher, 2006; Williams et al., 2006), and is based on a number of working meetings with members of the SEEK project as well as participants from the broader ecological community. In Section 3, we outline applications

Site	Trans.	Sp.	Dist.	Ht.	Area	Crabs
bird	1	hya	0.55	0.46	1260	2
bird	1	hya	5.90	0.44	2830	4
bird	1	hya	19.35	0.13	180	0
bird	2	hya	3.55	0.32	5030	6
bird	2	hya	18.20	0.58	17670	18
bird	2	hya	29.75	0.15	310	1
south	1	hya	4.15	0.08	20	0
south	1	hya	15.00	0.34	4420	7
south	1	hya	20.05	0.45	6360	7
south	2	hya	12.30	0.40	1260	2
...

Species	Wth.	Dens.
A. hyacinthus	30	6
A. hyacinthus	40	19
A. hyacinthus	20	4
A. hyacinthus	20	4
A. gemmifera	10	3
A. gemmifera	20	1
A. gemmifera	20	2
A. gemmifera	10	3
A. gemmifera	10	1
A. palifera	10	0
...

Date	Site	Transect	Species	Distance	Length	Area	Crabs	Density.
12-08-2006	Bird Island	1	A. hyacinthus	0.55	0.5	1260	2	16
12-08-2006	Bird Island	1	A. hyacinthus	5.90	0.4	2830	4	14
12-08-2006	Bird Island	1	A. hyacinthus	19.35	0.1	180	1	0
12-08-2006	Bird Island	2	A. hyacinthus	3.55	0.3	5030	6	12
12-08-2006	Bird Island	2	A. hyacinthus	18.20	0.6	17670	13	10
12-08-2006	Bird Island	2	A. hyacinthus	29.75	0.2	310	4	32
12-08-2006	South Island	1	A. hyacinthus	4.15	0.1	20	3	0
12-08-2006	South Island	1	A. hyacinthus	15.00	0.3	4420	7	16
12-08-2006	South Island	1	A. hyacinthus	20.05	0.5	6360	7	11
12-08-2006	South Island	2	A. hyacinthus	12.30	0.4	1260	6	16
...
10-18-2006	South Island		A. hyacinthus		0.3			6
10-18-2006	South Island		A. hyacinthus		0.4			19
10-18-2006	South Island		A. hyacinthus		0.2			4
10-18-2006	South Island		A. hyacinthus		0.2			4
10-18-2006	South Island		A. gemmifera		0.1			3
10-18-2006	South Island		A. gemmifera		0.2			1
10-18-2006	South Island		A. gemmifera		0.2			2
10-18-2006	South Island		A. gemmifera		0.1			3
10-18-2006	South Island		A. gemmifera		0.1			1
10-18-2006	South Island		A. palifera		0.1			0
...

Fig. 2 – Example data sets where: I contains data about coral crabs living within different coral species, including the given location and replicate transect, the species name, distance along the transect, colony area for coral colonies, and the number of crabs found in each colony; II contains data (from another study) about the density of coral crabs in different coral colonies, including the species name and crab density (other attributes, such as date and location, are described in the metadata, e.g., field notes); and III contains the results of merging I and II.

of our ontology for data discovery and integration. Finally, Section 4 concludes with a summary of our contributions and future work.

2. The observation ontology

The goal of the *Extensible Observation Ontology* (OBOE) is to serve as a formal and generic conceptual framework for describing the semantics of observational data sets (i.e., data sets consisting of observations and measurements). OBOE also prescribes a *structured* approach for organizing domain-specific ontologies through the use of “extension points.” OBOE extension points allow ontology classes, properties, and constraints to be easily defined for a particular domain-specific terminology, and existing domain extensions to be interrelated. Thus, OBOE can serve as a *framework* for defining new domain ontologies as well as interoperating and relating existing ones. Fig. 1 graphically depicts the basic core structure of OBOE, which consists of five classes labeled *Observation*, *Entity*, *Measurement*, *Characteristic*, and *Measurement Standard*, and six properties labeled *hasContext*, *ofEntity*, *hasMeasure-*

ment, *hasValue*, *hasPrecision*, *usesStandard*, and *ofCharacteristic*. Additional properties may be added to the *Measurement* class to capture when and where measurements were recorded, who recorded each measurement, the protocols of measurements, and so on. Similar properties can also be added to the *Observation* class. (Note that we capitalize OBOE classes to distinguish them from more general concepts, e.g., ‘Observation’ denotes an OBOE class whereas ‘observation’ denotes the more general concept.)

Today, most details of observational data are not recorded. Instead, the physical representation of data is often optimized for data collection; for use within a specific tool, e.g., R (*R Development Core Team, 2005*) or SAS® Software; or for a particular analysis, e.g., to perform a calculation requiring a site-by-species matrix. As a consequence, contextual information concerning data is typically *implicit*, where context is (possibly) encoded by attribute labels, implied by the proximity of attributes (i.e., neighboring data), stored in metadata as natural-language descriptions, or altogether missing. Consider the first data table in Fig. 2. The column of data labeled “Ht” can be assumed to represent height, giving information about the *characteristic* of some entity that was measured. However,

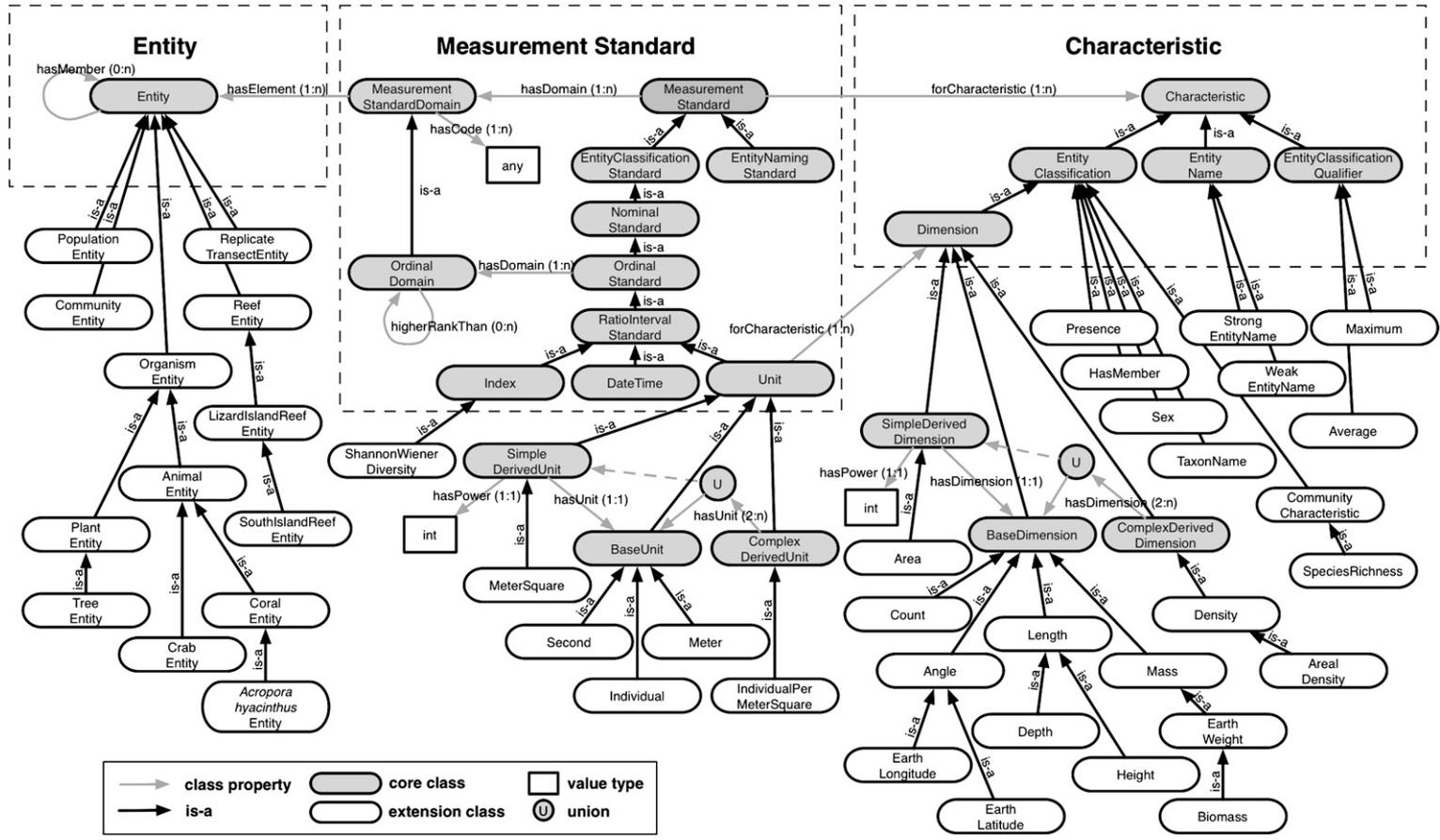


Fig. 3 – Detailed representation of OBOE core classes and examples of ontology extensions for Entity, Characteristic, and Measurement Standard. Shaded ellipses represent OBOE core classes, and open ellipses represent domain extensions for the examples given in this paper.

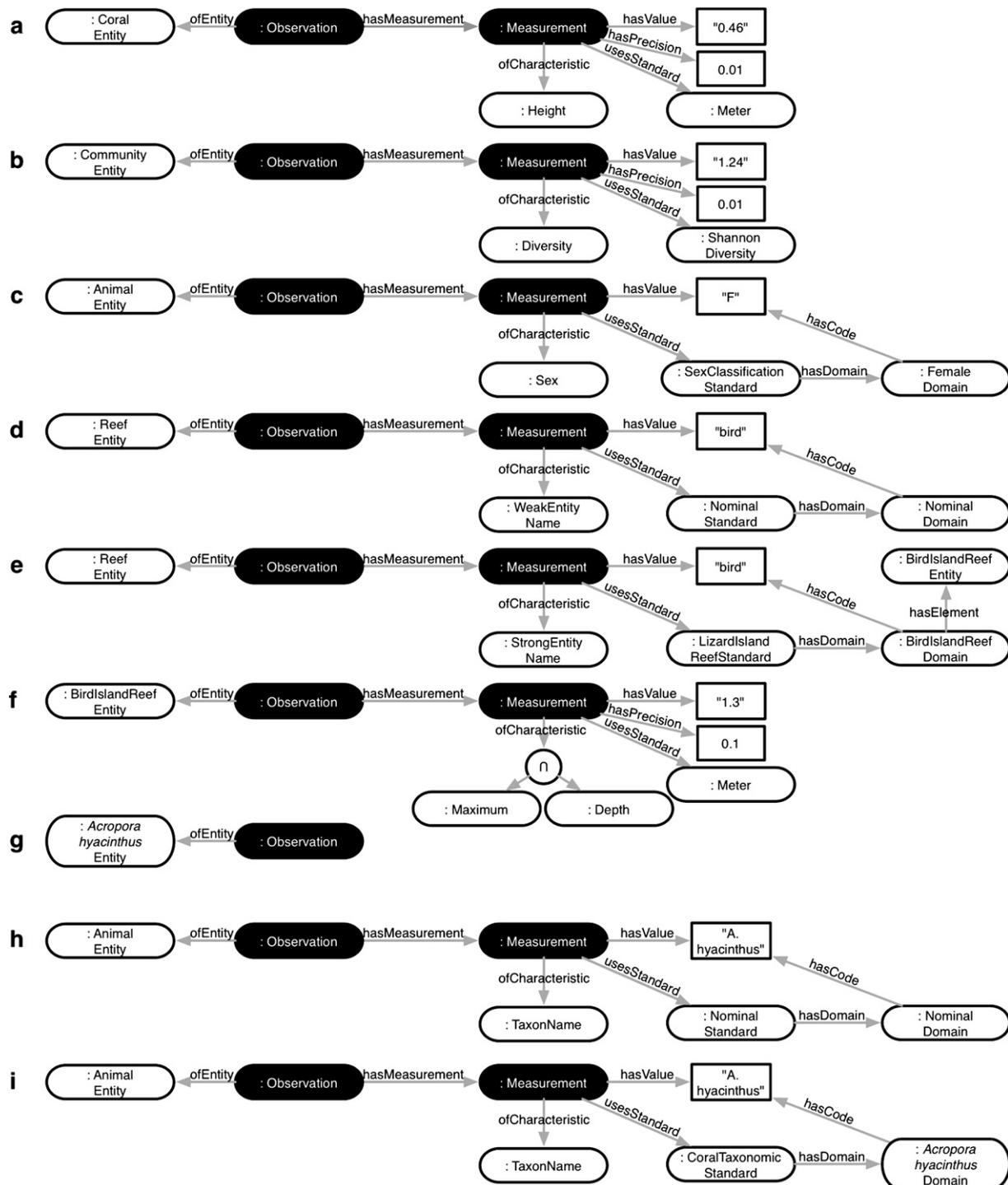


Fig. 4 – Examples of OBOE observational-data instantiations: (a) represents the observation of a coral where the coral height is measured in meters; (b) represents the observation of a community (an ecological concept) where diversity of taxa is measured using the Shannon Diversity index; (c) represents the observation of an animal classified as female via reference to a classification standard; (d) represents the observation of a reef entity classified as a “bird” (non-identifying), without reference to a specific Entity Classification Standard; (e) represents an observation of a reef with the (identifying) name of a particular reef according to an Entity Naming Standard for Lizard Island reefs; (f) represents the maximum depth of a particular reef; (g) represents an observation of the *Acropora hyacinthus* entity (i.e., the species concept) for which measurements can potentially be recorded; (h) represents an animal classified as a taxon, but without reference to an Entity Classification Standard, whereas for (i) a Standard is given.

no explicit information is given about the entity itself. The neighboring column labeled “Sp” suggests that each height measurement was for a species given in the column. Further, if all such species values in the column correspond to types of reef coral, one could surmise the kind of entities to which the height measurements pertain. The goal of OBOE is to provide a generic model for making such information explicit, which can then enable automated approaches for merging data (e.g., in this case with other coral data) and data discovery (e.g., for researchers looking for data on coral heights).

The rest of this section details the core OBOE classes and properties shown in Fig. 1, and demonstrates the use of OBOE for capturing observation data and for extending OBOE with domain ontologies. Although not discussed here, OBOE is encoded using the OWL-DL ontology language (McGuinness and van Harmelen, 2004), which gives a description-logic (Baader et al., 2003) formalization of the various parts of the OBOE ontology described below.

2.1. Observations

In OBOE, an observation is a statement that an entity of a particular type was observed. As shown in Fig. 1, all Observations are composed of exactly one Entity (expressed by the cardinality restriction “1:1” on the *ofEntity* property). The Entity class in OBOE represents all concrete and conceptual objects that are “observable.” While this notion of entity is extremely generic, it serves as a placeholder (i.e., extension point) for more specific types of objects. The left-hand portion of Fig. 3 gives a simple example of an Entity-class extension model, which is used in the examples of this paper. As shown, Entity classes are extended via *is-a* relations. For example, every *Organism Entity* is also an Entity, every *Plant Entity* is an Organism Entity (and hence an Entity), and so on. Hierarchies, like the Entity class hierarchy of Fig. 3, can be additionally constrained in OBOE using OWL-DL language constructs, specifying that classes are disjoint, equivalent, or related to multiple other classes combined through set union and intersection operations (e.g., stating that one class is equivalent to the union of two or more other classes). Although not shown, a number of Entity classes in Fig. 3 are defined as non-overlapping (disjoint). For instance, the *Plant* and *Animal Entity* classes are represented as disjoint sets of objects implying that no Plant Entities are Animal Entities, and vice versa. While the domain extensions portrayed in Fig. 3 are narrow and specialized for the purposes of this paper, we do not anticipate a single, universally accepted OBOE Entity classification. Instead, we aim to support multiple domain-specific extensions through OBOE, in which scientific groups and communities can flexibly build, share, and extend their own specialized entity (and other extension) models.

The Entity classes associated with observations in the first example data set of Fig. 2 include *Reefs*, *Replicate Transects*, *Animals*, and *Populations* of coral crab. In addition, some of these observations serve as context for other observations. For instance, each observation of a coral crab population occurs within the context of an animal (i.e., a coral colony). The *hasContext* property shown in Fig. 1 is used to capture these kinds of contextual relationships. In particular, an observation can serve as context for zero or more other observations, and can

itself have multiple contexts (e.g., a replicate transect may occur within a broader spatial context as well as a particular temporal context). Context in OBOE is defined independently from the observed entity, allowing the notion of observation scale to be efficiently formalized as well as systems implementing OBOE to perform automatic re-contextualization when merging observations (e.g., see Villa, 2007). Examples using context for merging observations are shown in Section 3.3.

The *hasContext* property asserts a “dependency” relationship between corresponding entities at the time of the observation, and thus, is defined to be transitive (Smith, 1996). For instance, each observation serving as context for replicate transect observations in the first example data set of Fig. 2 is also context for corresponding coral-colony observations. The *hasContext* property can also be extended to represent more specific types of contextualization, e.g., many of the mereological (i.e., part-of or containment) relations defined in Wand et al. (1999) would be suitable extensions. The transitive nature of *hasContext* can simplify the semantic annotation process in that by specifying only direct observation dependencies, it is possible to automatically infer all other indirect dependencies.

In OBOE, the type given to an observed entity is considered an essential quality (Guarino and Welty, 2002). An entity’s essential qualities help to define it, and always hold (are invariant) regardless of the entity’s context. For example, observing an entity of type *Animal Entity* implies that regardless of context, the particular object being observed is an “animal”. In this case, if the object was not an animal, it would be a different object altogether. Alternatively, one may assert that for a given context, a particular object is observed to be of a certain type, possibly to some degree of confidence. In this case, the type may not be an essential quality of the entity, since in a different context the quality may not hold. For instance, an object’s observed height may change in different contexts. Assuming height is not an essential quality of the object, it would not be correct to assign it a type such as *Tall Animal Entity*. Qualities that are not essential are attributed to entities in OBOE through measurements, which we describe in the following subsection.

2.2. Measurements

In OBOE, Observations can be composed of *Measurements*, which represent measurable *Characteristics* (i.e., qualities) of the entity being observed. As mentioned above, measurements in OBOE are assertions about an entity, and are not necessarily essential to the entity. Although not shown in Fig. 1, a Measurement is always associated to an Observation (i.e., no Measurement can exist without an associated Observation). Moreover, a particular Measurement can be associated with at most one Observation (i.e., two Observations may not share the same Measurement). Measurements assign values, via a *Measurement Standard*, to the characteristic of the associated entity. For certain types of measurements (e.g., physical quantities), a *Precision* is also given. Properties also exist for who recorded the measurement, when and where, and using what protocol (shaded properties, Fig. 1). These properties are not considered context, but

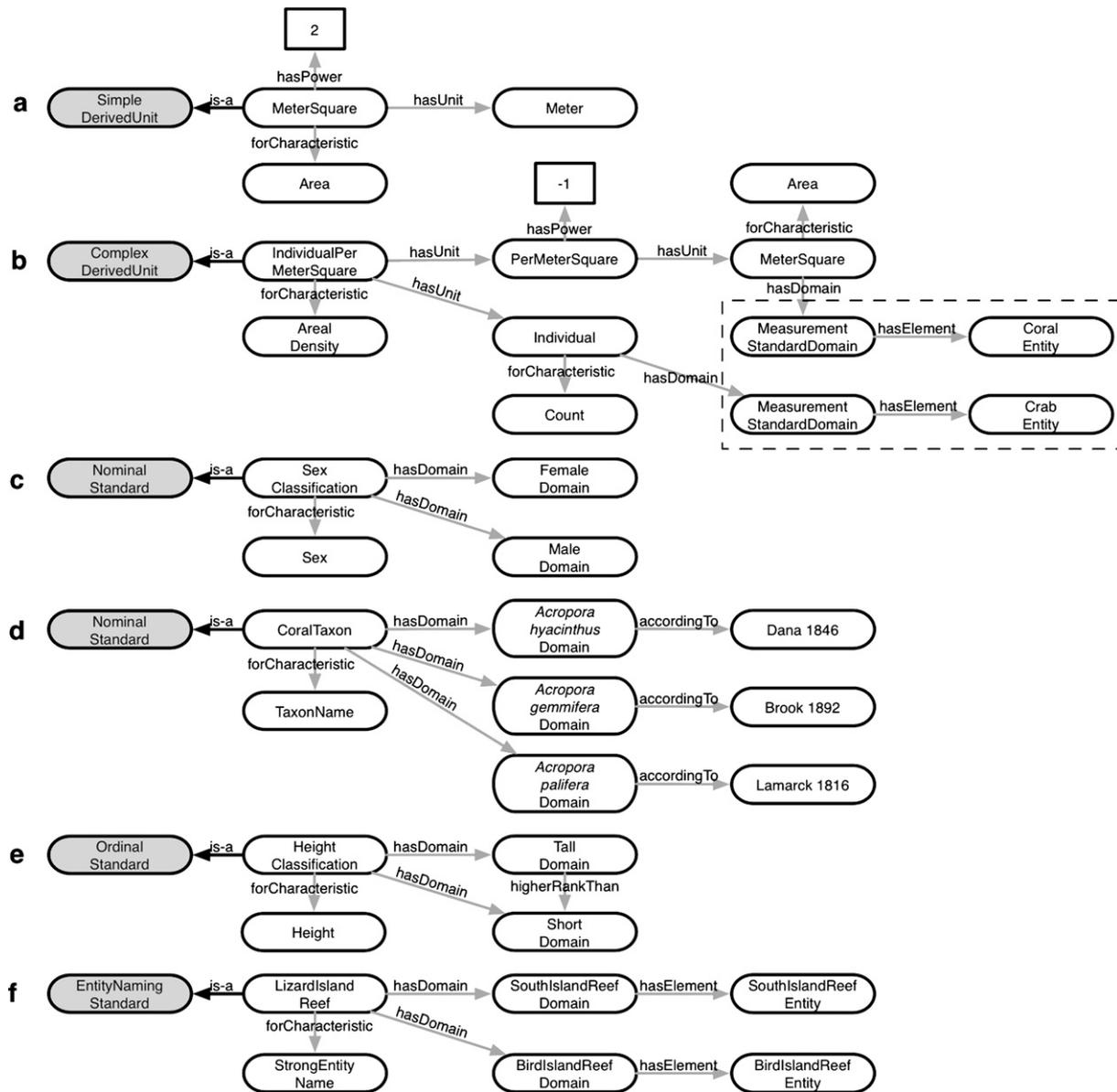


Fig. 5 – Example OBOE Measurement Standard extensions: (a) a Simple Derived Unit; (b) a Complex Derived Unit composed of two other units, as well as corresponding “semantic units” (dashed box), i.e., units linked to particular types of entities; (c and d) Nominal Classification Standards, one for male/female domains and another for taxonomic classification with “according to” relations; (e) an Ordinal Classification Standard including the “higherRankThan” property; and (f) an Entity Naming Standard representing a unique reef entities.

rather measurement-level metadata that can potentially be used for data discovery (e.g., finding all data that were measured using a particular protocol). Our definition of context captures observation-level transient hierarchies, which include observations of space and time for analytical purposes (e.g., the date was recorded because it was considered an important factor in an ecological study, but who recorded the date and when and where the date was recorded is not important from an analytical perspective).

Measurements in OBOE can be used to denote a wide range of instantiations of entity characteristics, including name or identity, a classification such as color (i.e., nominal and ordinal measurements; Stevens, 1946) or existence (e.g.,

measured as presence). Thus, the concept of measurement in OBOE is more generic than physical measurements alone. Fig. 4 gives a number of example observation-measurement instantiations. Note that we use the UML convention “id : Observation”, or simply “: Observation” when the id is unknown, to denote an instance of the Observation class (Jacobson et al., 1998). Fig. 4a states that a coral was observed such that the height of the coral was measured as 0.46 m with a precision of 0.01. Fig. 4b states that the diversity (a characteristic) of a community was measured as 1.24 according to the Shannon-diversity index.

In OBOE, Characteristics represent the types of measurable traits of entities, and denote another OBOE extension point.

Example Characteristic class extensions are shown at the right of Fig. 3, again defined via the *is-a* relation. *Dimension* is a core OBOE class that represents the distinguished set of physical dimensions that are used to define unit systems (described further below). Entity Classification and Name classes in our example extension ontology denote “stand-alone” trait types, e.g., which can be used directly in semantic annotations. Entity Classification Qualifier classes, however, must be combined—through intersection—with Entity Classification characteristics to be used in annotation. For instance, the class *Maximum Depth* used in Fig. 4f is defined as the intersection of the *Maximum* and *Depth* Characteristic classes. Unlike Entity Classification Characteristics, which group entities according to shared values (e.g., the entities of a particular depth), the Entity Name Characteristics are used to uniquely identify individual entities. *Weak Entity Names* uniquely identify entities that share observation context, i.e., weak names are context dependent. *Strong Entity Names* uniquely identify entities regardless of observation context. Plot names (e.g., Plot “A”, Plot “B”, etc.) are typical examples of weak entity names, and geographic locations are typical examples of strong entity names.

Characteristics and Entities are assumed to be disjoint (similar to Parsons and Wand (2000) and Bunge (1977)), even though the same term can often be used to refer to a characteristic or entity in natural language. For instance, ‘area’ could be the subject (i.e., entity) of an observation, where characteristics of the area entity are measured (e.g., its width). Conversely, ‘area’ could be a characteristic measured about an entity (e.g., the physical size of a study plot). Unique names are required in OBOE to distinguish these terms, e.g., *Area Entity* and *Area Characteristic*.

2.3. Measurement standards

A measured value (or “data point” in an observational data set) cannot be interpreted without reference to a defined measurement standard. Moreover, data integration relies on the ability to determine if two values are compatible, and if conversion to a common standard is possible. In OBOE, *Measurement Standards* are all the units, scales, categories, catalogs, and lists that are utilized when measuring a characteristic. Fig. 3 illustrates core OBOE classes for representing measurement standards, as well as a selection of classes used by examples in this paper (i.e., as a Measurement-Standard class extension). OBOE enforces a constraint (*forCharacteristic*) between the measurement standards and the characteristics that they represent. For example, the physical unit meter can only be used to represent characteristics that belong to dimension length, such as height.

OBOE defines two subclasses of Measurement Standard: Entity Naming Standard and Entity Classification Standard. An Entity Naming Standard is a naming scheme for globally identifying individual entities, where each entity is assumed to have only one instance (e.g., California or ID6547); whereas an Entity Classification Standard is for classifying entities by their traits, where each entity associated with a classification is assumed to have one or more instances (e.g., Tall, 12 m, Nitrogen Treatment, or Red). The Entity Classification Stan-

dard class contains the Steven’s Scale hierarchy: Nominal Standard (classifications are either the same or different), Ordinal Standard (classifications can be greater than or less than), and RatioInterval Standard (classifications are quantitative). The RatioInterval Standard is further subdivided into Unit, DateTime, and Index classes. The first of these, *Unit*, is subdivided into three disjoint classes. *Base Unit* contains the fundamental physical units, including SI units (e.g., meter, kilogram, second), and all manifestations of these units (e.g., millimeter, gram, hour). *Base Unit* also contains units for angle (e.g., degree), as well as number of items (individual). *Simple Derived Unit* contains all the physical units that are raised to a power other than 1, e.g., *Meter Square* (see Fig. 5a). The final unit class, *Complex Derived Unit*, is composed of two or more *Simple Derived Units* and/or *Base Units*, an example of which is given in Fig. 5b. Here, the *Individual per Meter Square* class is composed of the *Individual Base-Unit* class and the *Per Meter Square Simple Derived-Unit* class (i.e., *Meter Square* raised to the power of -1). The OBOE representation for units as well as the corresponding representation for dimensions (given in Fig. 3) is adopted from the approach used by the EML unit dictionary (Jones et al., 2001; Michener et al., 1997), which is based on the STXML language (Murray-Rust and Rzepa, 2002) and the NIST Reference on Constants, Units, and Uncertainty (<http://physics.nist.gov/cuu/Units/introduction.html>).

Measurement standards are composed of one or more measurement-standard domains through the *hasDomain* property. Each domain represents a possible value of the standard, and is (implicitly) related to the set of entities that have the corresponding characteristic value through the *hasElement* property. Measurement standard domains can be used to restrict the type of entities being measured. In the case of complex derived units (e.g., *areal density*, Fig. 5b)—composed from two or more independent unit types—each independent unit type can additionally be related to a corresponding entity type. In this way, OBOE provides a mechanism to describe so-called *semantic units*, e.g., grams of carbon per cubic meter of seawater, or individuals of rabbit per individuals of fox. For instance, in Table II of Fig. 2, the data value 6 represents the density of individuals in a coral, which can be expressed using a semantic unit (corresponding to the *Individual per Meter Square* class) defined as the composition of two unit components: the unit component *Individual* linked to the *Crab Entity* class, and the component *Meter Square* (raised to the power of -1) linked to the *Coral Entity* class (shown as dashed ellipses in Fig. 5b). Semantic units are commonly used in ecological data sets, and using OBOE it is possible to formally define the meaning of these units, making them available for discovery and integration processes (see Section 3).

The second subclass of Measurement Standard, *Index*, is a container for all the indices, scales, and surrogates of non-dimensional measures such as pH, the Richter Scale, and the various representations of biological diversity and evenness. *Indices* are often calculated using physical units, but have lost physical dimensionality due to a functional transformation (e.g., logarithmic or exponential transformation). For example, by log-transforming a measurement of height, dimensionality is lost, and the measurement instead becomes an index for height. Note that we do not consider a “dimensionless” unit of measurement in OBOE, since units are essential information

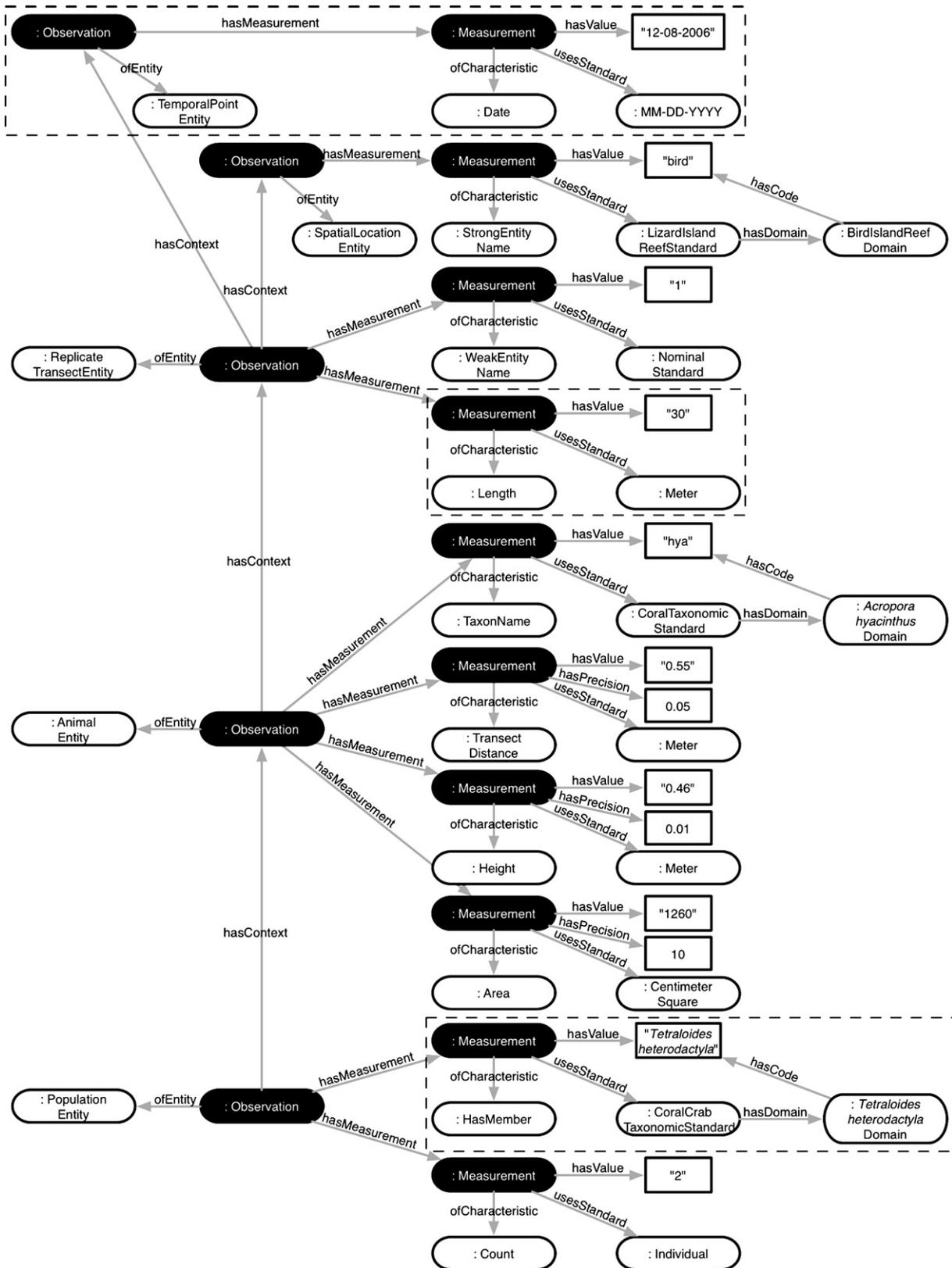


Fig. 6 – OBOE representation of the first row of data in Table I in Fig. 2. The colon placed before a class name denotes the instantiation of the class, i.e., “: Distance” is a member of the ontology class “Distance.” Solid ellipses represent the observation and measurement structure of the data set, including metadata (dashed box). Open ellipses represent terms selected from ontology extensions for each of Entity, Characteristics, and Measurement Standard. White squares represent the physical values from the data table or the metadata, and Precision where applicable.

in data integration even when the components of a complex unit cancel. For example, a ratio of two heights measured in meters is represented using a complex derived unit consisting of the base unit *meter* and the simple derived unit that represents *meter* raised to the power of -1 .

The third subclass, *DateTime*, is a container for the different ways of expressing a point in time (e.g., 12:30 am, 12-23-06, or 234 mya [million years ago]), as opposed to time intervals, which are represented using physical units (e.g., second, year, etc.). Many standards for representing the notion of date-time exist, and can be used to extend OBOE's *DateTime* class. Extensions for the *DateTime* class can include standards for representing Gregorian and Julian calendars, geological time scales, and the correlations among them. Measurements that use the first three subclasses of *RatioInterval Standard—Unit, Index, and DateTime*—also have *precision*, which indicates an estimate of the proximity of the measurement value to the real world value. Precision differs from accuracy, which is a methodological consideration, and thus, like methods for observation and measurement, is not elaborated here.

The subclass *Entity Classification Standard* contains measurement standards representing holistic characteristics of entities, such as taxonomic classification, color, or rank. Domains of classification standards are used for comparison when making a measurement (be it a conscious or subconscious comparison) and are assigned codes via the *hasCode* property denoting values of the standard. For example, the *Sex Classification Standard* class for the *Sex* characteristic class has the two domains *Female Domain* and *Male Domain*. Fig. 4c illustrates the use of this classification standard, stating that an animal was observed, the sex characteristic was measured, resulting in the value “F”, which is the code used in the data set for a *Female Entity* (from the classification domain). Therefore, this example asserts that the entity, in the particular observation context, is the intersection of animal (given as an essential quality) and female (asserted via measurement).

OBOE requires a measurement standard for all measurements. However, in certain situations a measurement standard does not exist, e.g., when an *ad hoc* naming scheme is used in a data set. The default OBOE measurement standard is *Entity Classification Standard* which assumes a corresponding *Classification Domain* whose code is the value of the measurement. For instance, Fig. 4d states that a particular reef was observed and that its name (in the current context) is “bird”, however, a reference to a known standard is not given. In this case, a default standard is used containing a single *Classification Domain* with the code “bird”. Here, the observed entity is not uniquely identified by the name “bird” since it is possible that many entities share this name; whereas the measurement given in Fig. 4e uses the *Lizard Island Reef* naming standard, which associates exactly one entity to the name “bird”, to uniquely identify the reef entity.

Finally, Fig. 4g–i also illustrates various representations of taxonomic name using OBOE. Fig. 4g considers *Acropora hyacinthus* the subject of observation (and therefore an essential quality), e.g., for describing measurements about a taxonomic concept or specimen. Fig. 4h considers the case where a relevant taxonomic domain is absent. Fig. 4i demonstrates the use of prescribing additional attributes to

an entity, in which data pertaining to *Acropora hyacinthus* is described as being of the taxonomic concept denoted by the *A. hyacinthus* Domain according to “Dana 1846” in the *Coral Taxonomic Standard* (Fig. 5d); where *taxonomic concept* denotes the representation standard defined in Kennedy et al. (2005).

3. Applications of the observation ontology

This section gives an overview of ways that OBOE can be used to facilitate the discovery and integration of ecological data. To expose the semantics of data using OBOE, data must first be annotated with relevant terms and relationships from the ontology. Semantic annotation is the process by which data are mapped to the ontology, and is accomplished by asserting the membership of data in ontology classes (e.g., *Animal Entity, Height, Meter*) and any additional relationships among classes (e.g., *hasContext*). An annotation language has been developed (e.g., see [Berkley et al., 2005; Bowers and Ludäscher, 2006]) that formalizes this mapping, and allows the annotation to be applied flexibly to multiple rows or cells depending on the data structure (e.g., table, cross-tabulation, etc.). Annotations expressed in this language are serialized according to an XML Schema, e.g., allowing them to be embedded within existing EML documents (or alternatively, as stand-alone documents that reference the corresponding EML). A graphical user interface for selecting different domain ontologies and annotating ecological data sets is being developed, which draw elements from EML to reduce the effort on the behalf of the user, and will eventually be integrated into a data set markup wizard.

Fig. 6 is a graphical representation of how the first row of the first data table (Table I in Fig. 2) might be annotated according to OBOE (Fig. 1) and the example domain extensions given in this paper (Fig. 3). The annotation in Fig. 6 would be applied to every row of the data in Table I. The left hand region of the figure shows the contextual hierarchy of observed entities. The observation of *Temporal Point* comes from the EML accompanying the raw data set (dashed box), and applies to the whole data set. This temporal observation, and a *Spatial Location* observation of the reef where the study took place that was recorded in the raw data, together provide independent context for the *Field Site* (i.e., the time the study took place is not dependent on the place, and visa versa). *Field Site* in turn provides context for the *Replicate Transect* (corresponding with “Site” and “Trans.” in Table I, respectively). As mentioned above, because context is transitive, *Temporal Point* also provides context for the replicate transect, although the observation of the transect still depends on that of the Reef. Further, the replicate transect provides context for the observation of coral colony, and the coral provides context for the observation of a crab population.

At each level of observation hierarchy, measurements were taken. Note that *measurements* tend to represent single columns in data tables (or else are derived from metadata), but *observations* tend to span one or more columns. For example, when the replicate transect was observed, only its name was measured to indicate that it differed from other replicate transects at a given reef. However, additional knowledge about the replicate transect was recorded, possibly

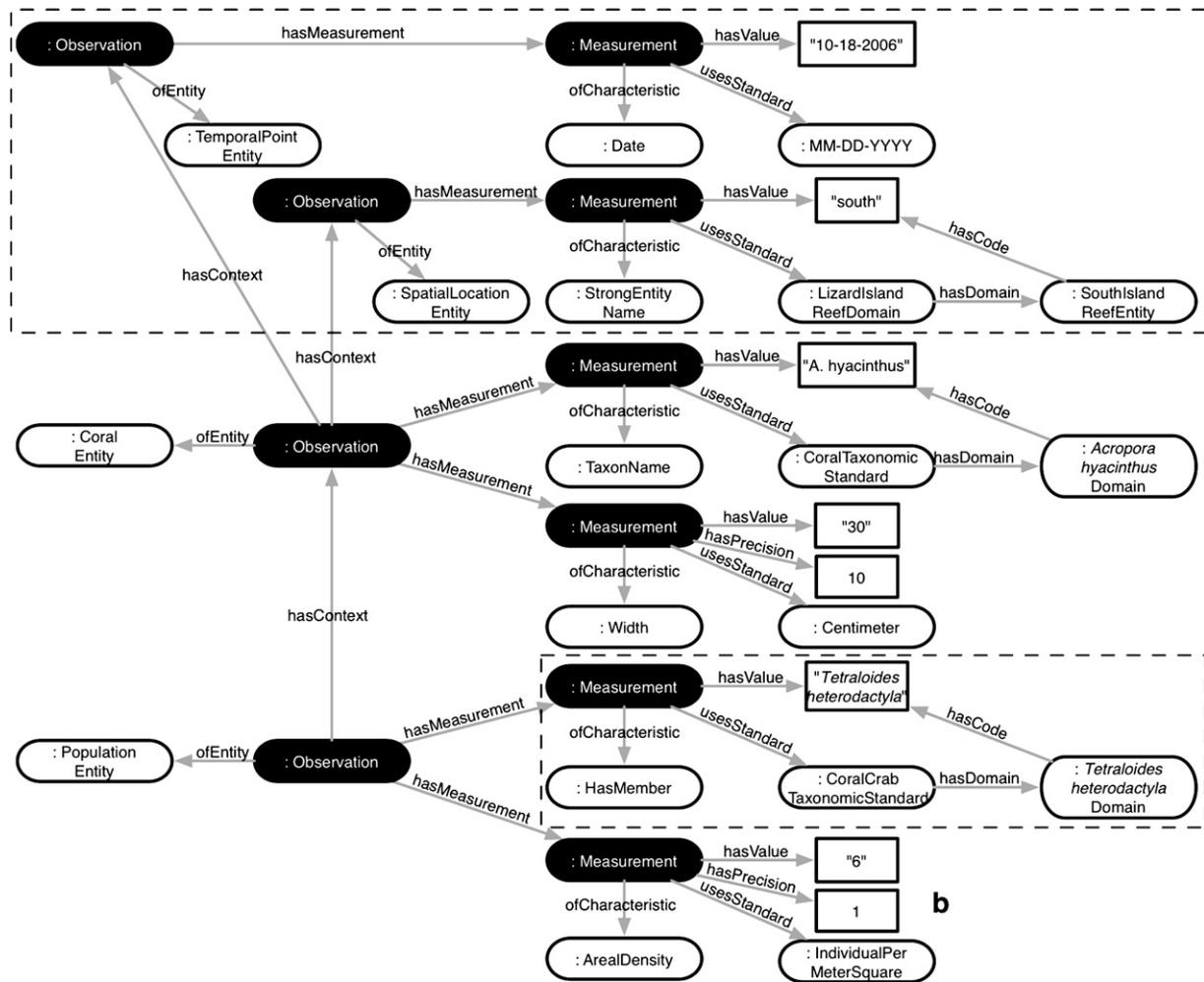


Fig. 7 – Formal OBOE representation of the first row of data in Table II in Fig. 2. Dashed box represents metadata. **b** corresponds with example **b** in Fig. 5, illustrating the more detailed representation using the semantic unit. See text for details.

in field notes, *e.g.*, that the transect was 30 m long. On the other hand, four measurements (representing four separate columns of data) were taken for a given coral entity found on the replicate transect: (1) the coral colony's taxonomic name according to a given taxonomic domain, (2) the distance along the transect where the colony was found, (3) a measure of colony area, and (4) colony height. Meanwhile, within the coral colony, the population of coral crabs was observed, and the number of individuals measured. Because all the crab populations were of the same crab species, taxonomic name was absent from the raw data, but was recorded as metadata.

Fig. 7 is a graphical representation of the first row of the second data table (Table II in Fig. 2). Here, much of the semantic information is extracted from the meta-data (dashed box), *i.e.*, when and where the observations were conducted, including geographic coordinates. Within this space-time context, a coral colony was observed, and its taxonomic name and height measured. In turn, the coral colony provides context for an observation of a population of coral crabs, where the number of individuals per unit area of the colony was measured, however data about how this calculation was made is missing and therefore implicit. Fig. 5b

illustrates how the crab population density can be represented using a semantic unit (where sub-units are for different entities), which can aid the data integration process outlined below.

Using the semantic annotations shown in Figs. 6 and 7, the rest of this section illustrates three useful applications leveraging the formal structure of OBOE: data discovery, summarization, and integration.

3.1. Data discovery

A major application facilitated by OBOE is the capability to discover data sets based on the concepts they represent (*i.e.*, their semantics), rather than just the labels and keywords that are used in traditional searches (Berkley et al., 2001). As mentioned previously, a data attribute labeled "Ht" in the first row of Table I is ambiguous, even to human interpretation. However, annotating this attribute with the Height (a Characteristic) from OBOE clarifies its meaning and relationship with other ontology terms (*i.e.*, via *is-a* relations, *part-of* relations, and other description-logic constraints). An enhanced keyword search for data about "height", *e.g.*, can leverage OBOE

definitions to discover the various data sets annotated with ontology terms related to Height, such as those assigned more specific terms like Body Height (Fig. 3). Search can also exploit relationships defined in OBOE Entity extension models, including the use of *part-of* relationships between classes. A search for “coral”, *e.g.*, could include entities that are *part-of* a coral colony, such as branch, tissue, skeleton, polyp, and so on. Note here that transitivity (*e.g.*, of coral parts), and other description-logic constraints (set intersection, etc.) can also be exploited to further enhance search.

More complex forms of inference can also be used, leveraging the logical structure of OBOE. For example, measurement dimensionality can be exploited to enhance keyword search, where a search for “density” data sets will not only return those annotated with Areal Density, but also those data sets that contain the Count and Area dimensions via appropriately contextualized observations (see Fig. 10). That is, this search would discover not only data having Areal Density attributes (Fig. 6), but also data sets having an attribute for Area and another attribute, functionally dependent via context (Fig. 5), for Count. In general, we are exploring the use of OBOE in this way for providing enhanced data-discovery query results, in which keyword queries given by scientists are: (1) expanded into their corresponding ontology classes, similar to traditional approaches based on formal terminologies (Voorhees, 1994; Moldovan and Mihalcea, 2000; Jarvelin et al., 2001); and (2) these ontology classes are compared (*via is-a, part-of, and so on*) to the explicit (*e.g.*, Height) and implicit (*e.g.*, where Count and Area imply Density) classes expressed in semantic annotations.

3.2. Data summarization

Upon finding a potentially relevant data set, an important aspect of the discovery process is to rapidly understand the content of the data set, *e.g.*, to determine whether the data is relevant for a particular analysis. An often-used approach for understanding data content is to aggregate (*i.e.*, summarize) attributes at various combinations of observation and measurement. The OBOE framework can be used to suggest appropriate data summarizations, and in so doing, also determine when a particular summarization is “sensible”. This notion of determining sensible summarizations exploits the basic structure of OBOE (Fig. 8a). For example, measurements can only be “sensibly” aggregated by other measurements that are of the same entity or measurements of entities providing observation context (*i.e.*, “higher” contextual entities). Therefore, it makes “sense” to summarize animal height by reef name, because the observation of reef provides context for the observation of animal (Fig. 8b). However, if water temperature were recorded at the contextual level of the study site, it would not make “sense” to summarize temperature by animal taxon name, which gives an arbitrary average temperature of the reefs where a taxon was measured, and which is dependent on the disparity in the number of animals measured at reefs. Furthermore, weak entity names are context dependent, and so replicate transect “1” at the study location “bird” is not the same as transect “1” at the study location “south”. It is therefore not “sensible” to summarize animal height or population count by replicate transects,

unless study location (a strong entity name) is also taken into account (Fig. 8c).

In general, the logical structure and constraints of OBOE can be used to test the usefulness of various statistical operations and modeling procedures. For example, when summarizing a nominal or categorical variable, the aggregate count is applicable, whereas continuous aggregations (*e.g.*, sum, average, maximum, standard deviation, etc.) are not. Moreover, types of modeling approaches can be recommended based on measurement types. For example, parametric linear models require continuous measures (*i.e.*, using a Ratio Interval Standard), whereas the inclusion of a categorical variable (*i.e.*, using Measurement Standards other than Ratio Interval) requires non-parametric model fitting approaches. By making the semantics of observation and measurement explicit in a formal structure such as OBOE, automated inferencing procedures (*i.e.*, “machine reasoning”) can be applied to enhance the various decision making processes used in exploring and modeling observational data sets (Gray et al., 1997).

3.3. Data integration

Another major application that OBOE facilitates is the capability to determine if two data sets can be either fully or partially merged once they are discovered (and vetted, *e.g.*, through summarization). To do so, a number of steps must be taken, starting at the lowest level semantic resolution (*i.e.*, at the level of single measurements and observations), and building up to the data set level (*e.g.*, aligning observation context). As a simple example, assume we are interested in merging data about heights following the discovery of the data tables I and II in Fig. 2. Fig. 9 illustrates the reasoning process involved in merging the first height observations in each data table, which is based on the semantic annotations using OBOE shown in Figs. 6 and 7. The first step is to determine if, and at what semantic resolution, the data are compatible. At the observation level, instances of the Animal and Coral entities are compatible at the semantic the resolution of Animal. Observation b loses semantic resolution when standardized for the merge (*i.e.*, the Animal entity class is the “lowest common denominator”). Similarly for measurement, instances of the Height and Width characteristics are compatible at the resolution of Length (*i.e.*, Length is the “common ancestor” of the Height and Width characteristics), where both instances in this case lose semantic resolution. Finally, the measurement standards used for the two measurements have the same dimension, and therefore they are compatible under OBOE’s model. Assuming merging lengths of coral in this way is desirable (*e.g.*, for a particular analysis), Fig. 9 illustrates the resulting merged data, where both observations are now of type Animal, both measurements are of characteristic Length, both measurement standards are in Meters, and the values are recalculated at the coarsest level of precision.

As a more sophisticated example, assume the result of a discovery search for data about “animal density” on “reefs” returns the two example data sets in Fig. 2. The data in Table I were discovered for the reasons discussed above. To integrate the data pertaining to areal densities, the OBOE structure can be used to determine if the data are compatible and then

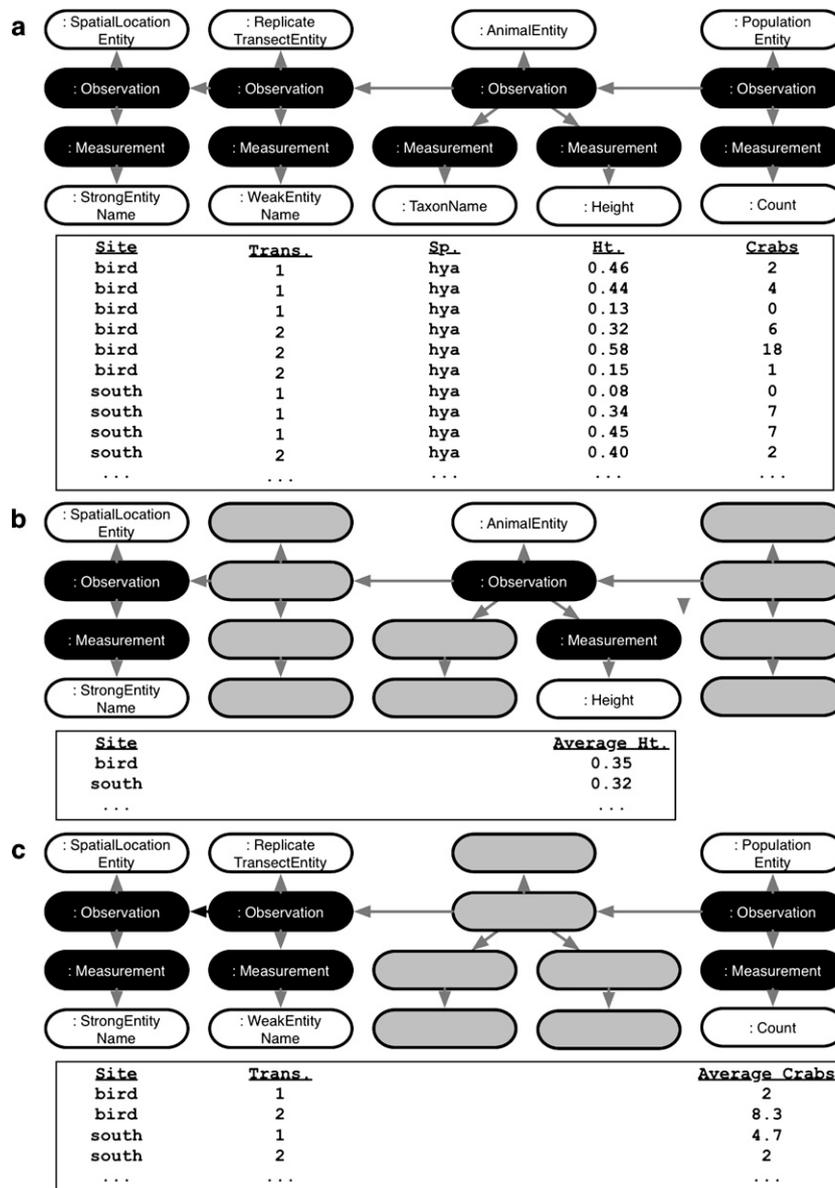


Fig. 8 – Graphical examples of data summarization leveraging OBOE's formal structure. (a) The basic observation and measurement structure for a subset of Table I. (b) Summarizing animal height by reef site ("sensible"). (c) Summarizing number of crab by animal species, replicate transect, and reef site (not "sensible"). (d) Summarizing average number of coral crabs by replicate transect name (not "sensible"). (e) Summarizing average number of coral crabs replicate transect ("sensible"). See text for details.

calculate the merge. The numerator in a semantic unit typically refers to the focal entity (count of individuals), whereas the denominator is contextual (area of coral). Therefore, semantic units are in essence a nested contextual dependency (required when more detailed information about context is missing). Therefore, if a focal entity is compatible with the numerator of the semantic unit, and the contextual entity is compatible with the denominator (e.g., Fig. 10a and b), the merge result shown at the bottom of Fig. 10 can be automatically computed. In the case of the areal densities of crab populations, Fig. 10a and b illustrate the semantic equivalencies between the two data tables, and calculates the merge, where a loses resolution when being coerced into the semantic unit form. However, knowledge about the coral

area and crab population in data Table I can still be retained (Table III), although there are no corresponding measurements from data Table II.

4. Summary

This paper presents the OBOE ontology framework for capturing the process of ecological field observation and measurement, facilitating logic-based reasoning (via description logic and OWL-DL) to be utilized to automate important data-management applications for data synthesis. The ontology formalizes an interpretation of observation, which focuses

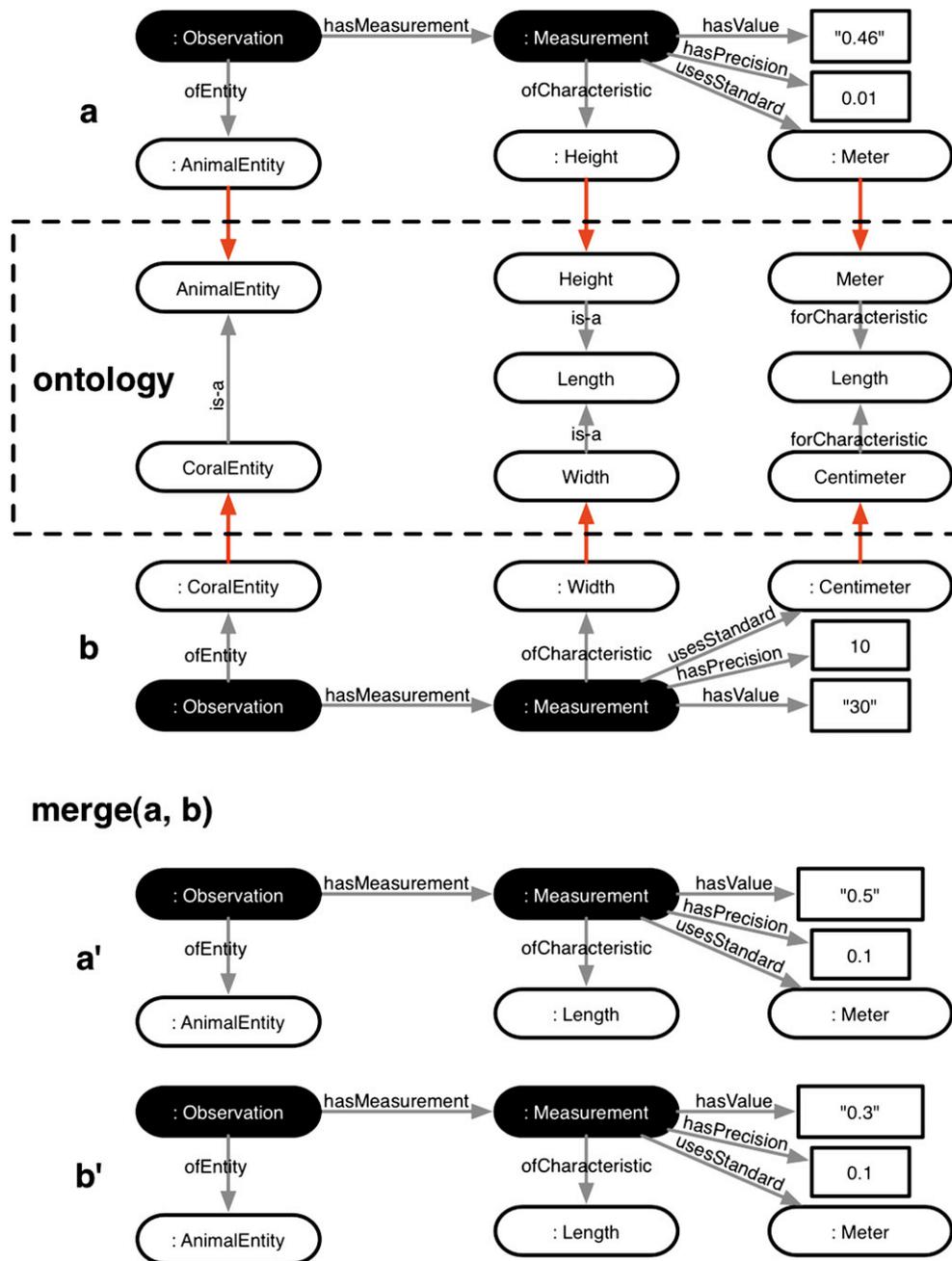


Fig. 9 – Merging measurements of animal height and coral length (from Tables I and II, respectively) based on semantic annotation with terms from OBOE (Figs. 6 and 7). See text for details.

on observing a concrete or conceptual entity, and measuring one or more of the entity’s characteristics by comparison with a measurement standard. In general, each data point in an observational data set is an instance of observation, and can provide contextual information for other instances of observation. Semantic annotations define a standard representation for mapping observational data to the ontology, and can be exploited in data discovery and integration applications. Annotation via OBOE and associated domain extensions makes explicit the basic definition of data and their relationships with other data, allowing annotated data sets to be easily contrasted. This approach can facilitate more powerful data

discovery and integration approaches, and can provide guidance for, and automate, data aggregation and summary.

Annotation to OBOE enables compatibility testing among data attributes, both at the level of the attribute (i.e., are the entity, characteristic, and measurement standard compatible?) and the data-set level (i.e., are entity nesting structures in two independent data sets compatible?). If compatible, the ontology contains the necessary details (i.e., constraints) to conduct the appropriate conversions so that data can be merged. Finally, OBOE provides a structured approach for creating domain-specific ontologies, allowing new ontologies to extend core OBOE classes. In this way, OBOE can also be used as a “glue structure”

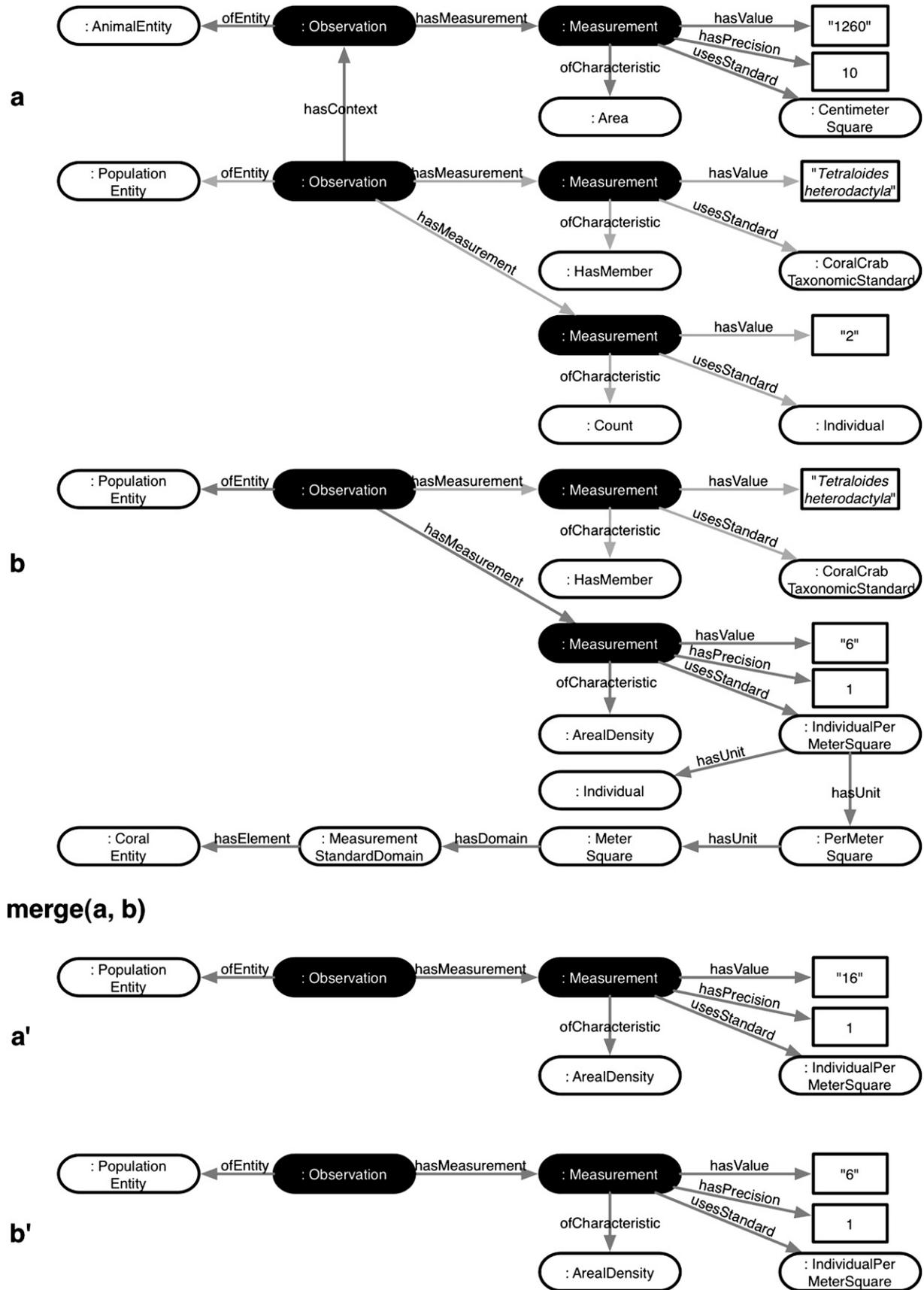


Fig. 10 – Merging measurements of crab population count and coral area (from Table I) and crab population areal density (from Tables II) based on semantic annotation with terms from OBOE (Figs. 6 and 7). See text for details.

to incorporate and inter-relate existing domain ontologies, allowing otherwise *ad hoc* ontologies to be structured and placed within a broader, cross-discipline scientific context.

Our ongoing and future work includes the development of an easy-to-use graphical user interface for data annotation based on OBOE. This tool will leverage existing EML metadata definitions (e.g., for basic data structure information and measurement units), will leverage OBOE ontology constraints to help direct and fill-in annotations when possible, and will transparently store semantic annotations using the XML serialization syntax mentioned in Section 3. We are also using OBOE as the foundation ontology in the SEEK Semantic Mediation System, which will take semantic annotations as input, and provide discovery, summarization, and integration services for use within, e.g., the Kepler scientific workflow system (Ludäscher et al., 2006; Berkley et al., 2005) and the EML-based Morpho application (Higgins et al., 2002).

REFERENCES

- Baader, F., Calvanese, D., McGuinness, D., Nardi, D., Patel-Schneider, P., 2003. *The Description Logic Handbook: Theory, Implementation, and Applications*. Cambridge University Press.
- Bard, J.R.L., Rhee, S.Y., 2004. Ontologies in biology: design, applications and future challenges. *Nat. Rev., Genet.* 5, 213–222.
- Batini, C., Ceri, S., Navathe, S.B., 1992. *Conceptual Database Design: An Entity-Relationship Approach*. Benjamin Cummings, Redwood City, CA.
- Berkley, C., Jones, M.B., Bojilova, J., Higgins, D., 2001. Metacat: a schema-independent XML database system. *Proc. of the 13th Intl. Conf. on Scientific and Statistical Database Management*. IEEE Computer Society.
- Berkley, C., Bowers, S., Jones, M.B., Ludäscher, B., Schildhauer, M., Tao, J., 2005. Incorporating semantics in scientific workflow authoring. *Proceedings of the 17th International Conference on Scientific and Statistical Database Management*. IEEE Computer Society.
- Bowers, S., Ludäscher, B., 2006. A calculus for propagating semantic annotations through scientific workflow queries. *Proc. of the Wkshp. on Query Languages and Query Processing*. *Lect. Notes Comput. Sci.*, vol. 4254, pp. 712–723.
- Bowers, S., Thau, D., Williams, R., Ludäscher, B., 2005. Data procurement for enabling scientific workflows: on exploring inter-ant parasitism. *Proc. of the Intl. Wkshp. on SemanticWeb and Databases*. *Lect. Notes Comput. Sci.*, vol. 3372.
- Brilhante, V.B., 2003. *Ontology and reuse in model synthesis*. PhD thesis. Univ. Edinburgh.
- Bunge, M., 1977. *Treatise on basic philosophy: vol. 1. Ontology 1: Furniture of the World*. Reidel, Boston.
- Chen, P.P., 1976. The entity-relationship model: toward a unified view of data. *ACM TODS* 1, 9–36.
- Cox, S. (Ed.), 2006. *Observations and Measurement*. Open GeoSpatial Consortium, Inc. OGC-05-087r4.
- Darwin Core, 2004. Darwin Core Schema (version 1.3), A draft standard of the Taxonomic Database Working Group (TDWG), <http://wiki.tdwg.org/DarwinCore>.
- DCMI, 2006. DCMI Metadata Terms. <http://www.dublincore.org/documents/dcmi-terms>.
- Ecological Metadata Language (EML) Specification. <http://knb.ecoinformatics.org/software/eml/>.
- Fegraus, E.H., Andelman, S., Jones, M.B., Schildhauer, M., 2005. Maximizing the value of ecological data with structured metadata: an introduction to ecological metadata language (EML) and Principles for metadata creation. *Bull. Ecol. Soc. Am.* 86, 158–168.
- Gray, J., Chaudhuri, S., Bosworth, A., Layman, A., Reichart, D., Venkatrao, M., Pellow, F., Pirahesh, H., 1997. Data cube: a relational aggregation operator generalizing group-by, cross-tab, and sub-totals. *Data Min. Knowl. Disc.* 1, 29–53.
- Gross, K.L., Pake, C.E., eds., (1995). *The Future of Long-term Ecological Data. A Report to the Ecological Society of America*. Vol. I: Text of the Report. 123 pp. Vol. II: Directories to Sources of Long-term Ecological Data. 114 pp.
- Gruber, T.R., Olsen, G.R., 1994. An ontology for engineering mathematics. *Proceedings of the Intl. Conf. on Principles of Knowledge Representation and Reasoning (KR)*, pp. 258–269.
- Guarino, N., Welty, C.A., 2002. Evaluating ontological decisions with OntoClean. *Commun. ACM* 45, 61–65.
- Hammer, J., McLeod, D., 1999. Resolution of representational diversity in multidatabase systems. In: Elmagarmid, A., Rusinkiewicz, M., Sheth, A. (Eds.), *Management of Heterogeneous and Autonomous Database Systems*, vol. 4. Morgan Kaufmann, San Francisco, pp. 91–117. 413 pp.
- Higgins, D., Berkley, C., Jones, M.B., 2002. Managing heterogeneous ecological data using morpho. *Proc. of the 14th Intl. Conf. on Scientific and Statistical Database Management*.
- Jacobson, I., Booch, G., Rumbaugh, J., 1998. *The Unified Software Development Process*. Addison Wesley Longman.
- Jarvelin, K., Kekalainen, J., Niemi, T., 2001. ExpansionTool: concept-based query expansion and construction. *Inf. Retr.* 4 (3–4), 231–255.
- Jones, M., Berkley, C., Bojilova, J., Schildhauer, M., 2001. Managing scientific metadata. *IEEE Internet Comput.* 5, 59–68.
- Jones, M.B., Schildhauer, M., Reichman, O.J., Bowers, S., 2006. The New Bioinformatics: integrating ecological data from the gene to the biosphere. *Ann. Rev. Ecol. Evol. Syst.* 37, 519–544.
- Keet, C.M., 2005. Factors affecting ontology development in ecology. *Lect. Notes Comput. Sci.* 3615, 46–62.
- Kennedy, J.B., Kukla, J., Paterson, T., 2005. Scientific names are ambiguous as identifiers for biological taxa: their context and definition are required for accurate data integration. *Lect. Notes Comput. Sci.* 3615, 80–95.
- Kepler Project, 2006. Kepler: An Extensible Scientific Workflow System. <http://kepler-project.org2006>.
- Ludäscher, B., Altintas, I., Berkley, C., Higgins, H., Jaeger, E., Jones, M., Lee, E.A., Tao, J., Zhao, Y., 2006. Scientific workflow management and the Kepler system. *Concurrency Comput. Pract. Exp.* 18 (10), 1039–1065.
- McGuinness, D.L., van Harmelen, F. eds., 2004. *OWL Web Ontology Language Overview*, W3C Recommendation 10 Feb 2004, <http://www.w3.org/TR/2004/REC-owl-features-20040210/>.
- Michener, W.K., 2000. Research design: translating ideas to data. In: Michener, W.K., Brunt, J.W. (Eds.), *Ecological Data: Design, Management and Processing*. Blackwell Science, Oxford, pp. 1–24.
- Michener, W.K., Brunt, J.W., Helly, J.J., Kirchner, T.B., Stafford, S.G., 1997. Non-geospatial metadata for the ecological sciences. *Ecol. Appl.* 7, 330–342.
- Moldovan, D.I., Mihalcea, R., 2000. Using WordNet and lexical operators to improve internet searchers. *IEEE Internet Comput.* 4 (1), 34–43.
- Murray-Rust, P., Rzepa, H.S., 2002. STXML: A markup language for scientific, technical and medical publishing. *Data Sci. J.* 1, 1–65.
- The Gene Ontology Consortium, 2000. Gene Ontology: tool for the unification of biology. *Nat. Rev., Genet.* 25, 25–29.
- The Observation Ontology, OBOE. <http://ecoinformatics.org/ontologies/observation-0.1.0>.
- Parsons, J., Wand, Y., 2000. Emancipating instances from the tyranny of classes in information modeling. *ACM Trans. Database Syst.* 25, 228–268.
- Pickett, S.T.A., Kolasa, J., Jones, C.G., 1994. *Ecological Understanding: The Nature of Theory and the Theory of Nature*. Academic Press, San Diego.

- Rosse, C., Mejino Jr., J.L.V., 2003. A reference ontology for bioinformatics: the Foundational Model of Anatomy. *J. Biomed. Inform.* 36, 478–500.
- R Development Core Team, 2005. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org>.
- Schentz, H., Mirtl, M., 2003. MORIS: a universal information system for environmental monitoring. In: Schimack, G.P. (Ed.), *Environment Software Systems*, vol. 5. Springer.
- Smith, B., 1996. Mereotopology: a theory of parts and boundaries. *Data Knowl. Eng.* 20, 287–303.
- Smith, B., 2001. Objects and their environments: from Aristotle to ecological ontology. In: Frank, A., Raper, J., Cheylan, J.P. (Eds.), *The Life and Motion of Socio-Economic Units*. Taylor and Francis, London, pp. 79–97.
- Smith, B., Varzi, A., 1999a. The niche. *Nous* 33 (2), 198–222.
- Smith, B., Varzi, A., 1999b. The formal structure of ecological contexts. In: Bouquet, P., Brezillon, P., Serafini, L., Beneceretti, M., Castellani, F. (Eds.), *Modeling and Using Context*. Lect. Notes Artif. Int., vol. 1688, pp. 339–350.
- Sowa, J.F., 1999. *Knowledge Representation: Logical, Philosophical, and Computational Foundations*. PWS Publishing Co., Boston.
- Stevens, S.S., 1946. On the theory of scales of measurement. *Science* 103, 677–680.
- Villa, F., 2007. A semantic framework and software design to enable the transparent integration, reorganization and discovery of natural systems knowledge. *J. Intell. Inf. Syst.* doi:10.1007/s10844-006-00322-x.
- Voorhees, E.M., 1994. Query expansion using lexical-semantic relations. *Proc. of the Annual Intl. ACM-SIGIR Conf. on Research and Development in Information Retrieval (SIGIR)*.
- Wand, Y., Storey, V.C., Weber, R., 1999. An ontological analysis of the relationship construct in conceptual modeling. *ACM Trans. Database Syst.* 24, 494–528.
- Williams, R.J., Martinez, N.D., Golbeck, J., 2006. Ontologies forecoinformatics. *J. Web Semant.* 4, 237–242.