

# Data Procurement for Enabling Scientific Workflows: On Exploring Inter-Ant Parasitism\*

Shawn Bowers<sup>1</sup>, David Thau<sup>2</sup>, Rich Williams<sup>3</sup>, and Bertram Ludäscher<sup>1</sup>

<sup>1</sup> San Diego Supercomputer Center, UCSD, La Jolla, CA, USA

<sup>2</sup> University of Kansas, Lawrence, KS, USA

<sup>3</sup> National Center for Ecological Analysis and Synthesis, UCSB, Santa Barbara, CA, USA

## 1 Introduction

Like content on the web, scientific data is highly heterogeneous and can benefit from rich semantic descriptions. In our work, we are particularly interested in developing an infrastructure for expressing explicit and fine-grain semantic descriptions of ecological data (and life-sciences data in general), and exploiting these descriptions to provide automated data integration and transformation within scientific workflows [2]. Using semantic descriptions, our goal is to provide scientists with: (1) tools to easily search for and retrieve datasets relevant to their study (i.e., data *procurement*), (2) the ability to select a subset of returned datasets as input to a scientific workflow, and (3) automated integration and restructuring of the selected datasets for seamless workflow execution.

As part of this effort, we are developing the *Semantic Mediation System* (SMS) within the SEEK project<sup>4</sup>, which aims at combining semantic-web technologies—namely OWL and RDF—with traditional data-integration techniques [3, 6, 7]. We observe that along with these “traditional” approaches, mediation of ecological data also requires external, special-purpose services for accessing information not easily or conveniently expressed using conceptual modeling languages, such as description logics. The following are two specific examples of ecologically relevant, external services that can be exploited for scientific-data integration and transformation.

**Taxonomic Classification and Mapping.** There is an extensive body of knowledge on species (both extinct and existing) represented in a variety of different taxonomic classifications, and new species are being discovered continually [9]. The same species can be denoted in many ways across different classifications, and resolving names of species requires mappings across multiple classification hierarchies [11]. Within SMS we want to leverage operations that exploit these existing mappings, e.g., to obtain synonyms of species names, without explicitly representing the mappings or simulating the associated operations within the mediator.

**Semantics-Based Data Conversion.** We are interested in applying operations during mediation that can transform and integrate data values based on their implied meaning. However, for scientific data, the nature of these conversions are often difficult to express explicitly within a conceptual model. A large number of ecological datasets represent real-world observations (like measuring the abundance of a particular species),

---

\* This work supported in part by NSF grant ITR 0225676.

<sup>4</sup> *Science Environment for Ecological Knowledge*, <http://seek.ecoinformatics.org>

and therefore often have slightly different spatial and temporal contexts, use different measurement protocols, and measure similar information in disparate ways (e.g., area and count in one dataset, and density, which is a function of area and count, in a second dataset). Like with taxonomic classification, we want the mediator to exploit existing conversion operations when possible.

We note that the application of semantics-based data conversion often depends on usage, i.e., some conversions may only be applied when certain conditions are met by the associated datasets. Thus, to correctly apply such conversions, SMS may require additional information to determine whether a particular conversion is applicable for a given dataset. In general, we believe that new techniques are required to support the use of these and similar external services within traditional data-integration architectures.

This short paper describes an initial logic-based SMS prototype that leverages ontologies, semantic descriptions, and simple external services (primarily taxonomic) to help researchers find relevant datasets for ecological modeling. The rest of this paper is organized as follows. In Section 2 we describe the motivating scenario for our SMS prototype. In Section 3 we discuss the details of the prototype. And finally, in Section 4 we conclude by discussing future work.

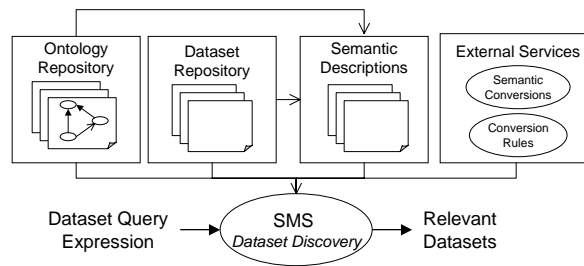
## 2 Motivation: Ant Parasitism and Niche Modeling

A diverse and much studied group of organisms in ecology is the family *Formicidae*, commonly known as ‘ants’. Ants generally account for ten to fifteen percent of the animal biomass of any given area. Beyond their important role in churning much of the earth’s soil, ants are social animals that provide insights into the evolution of social behaviors. One such complex social behavior is parasitism between ant species [4].

The environments in which parasitism is likely to occur provides important data on how parasitism arises. For example, one theory states that inter-ant parasitism is more likely to arise in colder climates than in warmer ones. Thus, an ecological researcher may be interested in testing the high-level question: *In California, based on existing data, which environmental properties play an important role in determining the ranges of ants involved in inter-ant parasitism?*

The verification of this question requires access to a wide array of data: (1) the types of parasitic relationships that exist between ants, (2) the names of species of ants taking part in these parasitic relationships, (3) georeferenced observations of these species of ants, and (4) the climate and other environmental data within the desired locations.

Today, these datasets are typically sought out by the researcher, retrieved, and integrated by hand. The researcher analyzes the data by running it through an appropriate ecological model, the result of which is used to help verify a hypothesis. In our example, an ecological niche model [8] can be used, which takes data about the presence of a species and the environmental conditions of the area in question, and produces a set of rules that define a “niche” (i.e., the conditions necessary for the species to exist) relative to the given environmental conditions and presence data. The rest of this paper describes a first step towards helping a researcher easily collect the datasets needed to test inter-ant parasitism, and similar high-level questions.



**Fig. 1.** The initial SMS architecture for ecological data mediation.

$d_1$	<u>genus</u>	<u>species</u>	<u>count</u>	<u>lat</u>	<u>lon</u>	$d_2$	<u>genus</u>	<u>species</u>	<u>cnt</u>	<u>lt</u>	<u>ln</u>
	Manica	parasitica	2	37.85	-119.57		Camponotus	forasini	1	-29.65	26.18
	Manica	bradelyi	1	38.32	-119.67						
$d_3$	<u>man-para-cnt</u>	<u>aph-cald-cnt</u>	<u>lt</u>	<u>ln</u>		$d_4$	<u>genus1</u>	<u>species1</u>	<u>genus2</u>	<u>species2</u>	
	3	6		37.56	-120.03		Manica	parasitica	Aphaenogaster	calderoni	

**Fig. 2.** Four heterogeneous datasets  $d_1$  through  $d_4$ .

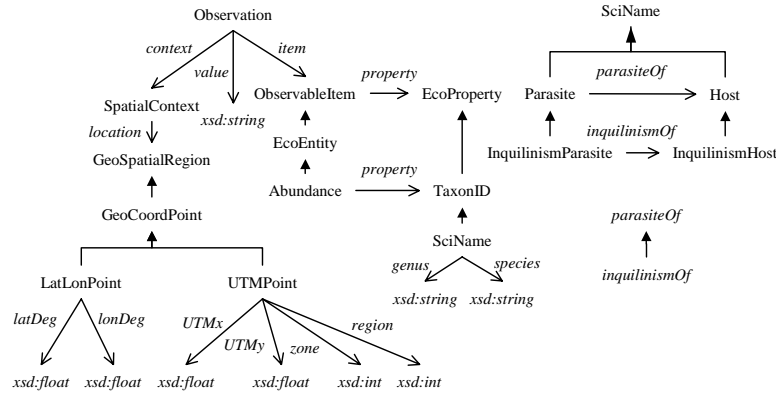
### 3 The Prototype

Our dataset-discovery architecture is shown in Figure 1. A set of repositories store ontological information, datasets, and semantic descriptions. Each dataset added to the repository is required to have a corresponding semantic description, expressed as a sound *local-as-view* mapping [6, 7] against concepts and roles in the ontologies. We also consider external services, which currently consist of synonym and unit-conversion operations. The SMS engine accepts a user query and returns the set of relevant datasets that satisfy the given query.

Figure 2 gives examples from four datasets used in analyses involving ant and inter-ant parasite presence data. Dataset  $d_1$  in Figure 2 contains georeferenced ant data from AntWeb<sup>5</sup> and consists of approximately seventeen-hundred observations, each of which consist of a genus and species scientific name, an abundance count, and the location of the observation. Dataset  $d_2$  in Figure 2 contains similar georeferenced ant data from the Iziko South African Museum (ISAM),<sup>6</sup> consisting of about twelve-thousand observations. Dataset  $d_3$  in Figure 2 is a typical representation used for georeferenced co-occurrence data, where species are encoded within the schema of the table. This dataset contains only five tuples. Dataset  $d_4$  in Figure 2 describes specific ants that participate in inquilinism inter-ant parasitism. The first two columns denote the parasite and the last two columns denote the host. Over two-hundred pairs of ants are described using four distinct datasets, each representing a particular parasitic relationship (all data were derived from Table 12-1 of [4]). Finally, Figure 3 shows a simplified fragment of the measurement and parasitism ontologies currently being developed within SEEK (where solid arrows denote *isa* relations).

<sup>5</sup> See [www.antweb.org](http://www.antweb.org)

<sup>6</sup> Provided by Hamish Robertson, Iziko Museums of Cape Town



**Fig. 3.** Simplified ontologies for measurement observations and inter-ant parasitism.

The following conjunctive queries define semantic descriptions of datasets  $d_1$ ,  $d_3$ , and  $d_4$  (note that the semantic description of  $d_2$  is identical to  $d_1$ ). Each semantic description expresses a *local-as-view* mapping [6, 7], defining a dataset in terms of the ontology of Figure 3.

$d_1(\text{Ge, Sp, Co, Lt, Ln}) :-$

Observation(O), value(O,Co), context(O,S), location(S,P), LatLonPoint(P),  
latDeg(P,Lt), lonDeg(P,Ln), item(O,A), Abundance(A), property(A,N), SciName(N),  
genus(N,Ge), species(N,Se).

$d_3(\text{Mp, Cf, Lt, Ln}) :-$

Observation( $O_1$ ), value( $O_1$ ,Mp), context( $O_1$ ,S), location(S,P), LatLonPoint(P),  
latDeg(P,Lt), lonDeg(P,Ln), item( $O_1$ , $A_1$ ), Abundance( $A_1$ ), property( $A_1$ , $N_1$ ),  
SciName( $N_1$ ), genus( $N_1$ , 'Manica'), species( $N_1$ , 'parasitica'), Observation( $O_2$ ),  
value( $O_2$ ,Cf), context( $O_2$ ,S), item( $O_2$ , $A_2$ ), Abundance( $A_2$ ), property( $A_2$ , $N_2$ ),  
SciName( $N_2$ ), genus( $N_2$ , 'Aphaenogaster'), species( $N_2$ , 'calderoni').

$d_4(\text{G}_1, \text{S}_1, \text{G}_2, \text{S}_2) :-$

InquilinismParasite(P), SciName(P), genus(P, $G_1$ ), species(P, $S_1$ ), InquilinismHost(H),  
genus(H, $G_2$ ), species(H, $S_2$ ), inquilinismOf(P,H).

The following example is a dataset-discovery query that finds all datasets containing georeferenced abundance measurements of *Manica bradleyi* ants observed within California (as defined by the given bounding box). Dataset-discovery queries allow predicates to be *annotated* with dataset variables, given as  $D$  below. A dataset handle is returned by the query if each formula annotated with  $D$  is satisfied by the dataset, assuming the given inequality (i.e., the latitude-longitude) conditions also hold.

$q_1(D) :-$  Observation(O) <sup>$D$</sup> , context(O,S) <sup>$D$</sup> , location(S,P) <sup>$D$</sup> , LatLonPoint(P) <sup>$D$</sup> ,  
latDeg(P,Lt) <sup>$D$</sup> , lonDeg(P,Ln) <sup>$D$</sup> , item(O,A) <sup>$D$</sup> , Abundance(A) <sup>$D$</sup> , property(A,N) <sup>$D$</sup> ,  
SciName(N) <sup>$D$</sup> , genus(N,'Manica') <sup>$D$</sup> , species(N,'bradleyi') <sup>$D$</sup> , Lt  $\geq$  33, Lt  $\leq$  42,  
Ln  $\geq$  -124.3, Ln  $\leq$  -115.

Using a standard data-integration query-answering algorithm [7], the query above is answered by (1) finding *relevant* information sources, i.e., sources whose view mappings

overlap with the given query, and (2) using the relevant sources, rewriting the user query, producing a sound query expressed only against the underlying data sources, possibly containing additional conditions. We extend this approach by also considering dataset annotations on query formulas. In our example,  $d_1$  and  $d_2$  are the only relevant datasets for the above query, giving the following query rewritings. Note that after executing the queries below, only  $d_1$  is returned; the ISAM dataset does not contain the given species.

$q_1(d_1) :- d_1(\text{'Manica', 'bradleyi', Ct, Lt, Ln}), Lt \geq 33, Lt \leq 42, Ln \geq -124.3, Ln \leq -115.$   
 $q_1(d_2) :- d_2(\text{'Manica', 'bradleyi', Ct, Lt, Ln}), Lt \geq 33, Lt \leq 42, Ln \geq -124.3, Ln \leq -115.$

The following query is similar to  $q_1$ , but uses an external service (prefixed with 'ext:') for computing synonymy of species names.

$q_2(D) :- \text{Observation}(O)^D, \text{context}(O,S)^D, \text{location}(S,P)^D, \text{LatLonPoint}(P)^D,$   
 $\text{latDeg}(P,Lt)^D, \text{lonDeg}(P,Ln)^D, \text{item}(O,A)^D, \text{Abundance}(A)^D, \text{property}(A,N)^D,$   
 $\text{SciName}(N)^D, \text{genus}(N,Ge)^D, \text{species}(N,Sp)^D, Lt \geq 33, Lt \leq 42, Ln \geq -124.3,$   
 $Ln \leq -115, \text{ext:synonym}(\text{'Manica', 'bradleyi', Ge, Sp}).$

The synonymy operation, encapsulated as a logical formula above, draws from descriptions (expressed as XML files) in the Hymenoptera Name Server [5], and supports over twenty-five hundred taxa of ants and their synonymy mappings. In the operation, a given genus-species pair is always a synonym of itself. We note that in the prototype, we equate synonyms between taxa as equivalence relations. This assumption is often an oversimplification [1] and in future work we intend to explore the impact of different synonymy relations between taxa.

The following rewritings are obtained from the above query. After execution, the rewritten  $q_2$  query will return dataset  $d_1$  as well as dataset  $d_3$ ; the latter because *Aphaenogaster calderoni* is a synonym of *Manica bradleyi*. Note that we could have discarded the third rewriting below since all arguments of the synonym operation are ground, and for the particular binding, the species' are not valid synonyms.

$q_2(d_1) :- d_1(\text{Ge, Sp, Ct, Lt, Ln}), Lt \geq 33, Lt \leq 42, Ln \geq -124.3, Ln \leq -115,$   
 $\text{ext:synonym}(\text{'Manica', 'bradleyi', Ge, Sp}).$   
 $q_2(d_2) :- d_2(\text{Ge, Sp, Ct, Lt, Ln}), Lt \geq 33, Lt \leq 42, Ln \geq -124.3, Ln \leq -115,$   
 $\text{ext:synonym}(\text{'Manica', 'bradleyi', Ge, Sp}).$   
 $q_2(d_3) :- d_3(\text{Mp, Cf, Lt, Ln}), Lt \geq 33, Lt \leq 42, Ln \geq -124.3, Ln \leq -115,$   
 $\text{ext:synonym}(\text{'Manica', 'bradleyi', 'Manica', 'parasitica'}).$   
 $q_2(d_3) :- d_3(\text{Mp, Cf, Lt, Ln}), Lt \geq 33, Lt \leq 42, Ln \geq -124.3, Ln \leq -115,$   
 $\text{ext:synonym}(\text{'Manica', 'bradleyi', 'Aphaenogaster', 'calderoni'}).$

Finally, the following query finds datasets containing georeferenced measurements of parasites of *Manica bradleyi* within California. Thus, the query finds the relevant ant presence data needed for our researcher's high-level question, for a single host species. The query uses the external synonym operation and projects the latitude, longitude, and genus and species names of the relevant observations (which could be used as "provenance" information, or for filtering the dataset for use in an analytical model).

$q_3(D,Lt,Ln,Ge,Sp) :-$  Observation(O)<sup>D</sup>, context(O,S)<sup>D</sup>, location(S,P)<sup>D</sup>, LatLonPoint(P)<sup>D</sup>, latDeg(P,Lt)<sup>D</sup>, lonDeg(P,Ln)<sup>D</sup>, item(O,A)<sup>D</sup>, Abundance(A)<sup>D</sup>, property(A,N)<sup>D</sup>, SciName(N)<sup>D</sup>, genus(N,Ge)<sup>D</sup>, species(N,Sp)<sup>D</sup>, Lt ≥ 32, Lt ≤ 42, Ln ≥ -124.3, Ln ≤ -115, Host(Ho), genus(Ho,Ge<sub>1</sub>,Sp<sub>1</sub>), ext:synonym('Manica','bradleyi',Ge<sub>1</sub>,Sp<sub>1</sub>), Parasite(Pa,Ge<sub>2</sub>,Sp<sub>2</sub>), parasiteOf(Pa,Ho), ext:synonym(Ge<sub>2</sub>,Sp<sub>2</sub>,Ge,Sp).

The rewritings of  $q_3$  are shown below. The result of executing the query will include the tuples (d<sub>1</sub>,37.85,-119.57,'Manica','parasitica') and (d<sub>3</sub>,37.56,-120.03,'Manica','parasitica'), where only datasets d<sub>1</sub> and d<sub>3</sub> will contain possible answers. In particular, Manica parasitica are inquilinism parasites of Manica bradleyi, which is derived from dataset d<sub>4</sub> by computing Manica bradleyi synonyms.

$q_3(d_1,Lt,Ln,Ge,Sp) :-$  d<sub>1</sub>(Ge,Sp,Ct,Lt,Ln), Lt ≥ 33, Lt ≤ 42, Ln ≥ -124.3, Ln ≤ -115, ext:synonym('Manica','bradleyi',Ge<sub>1</sub>,Sp<sub>1</sub>), d<sub>4</sub>(Ge<sub>1</sub>,Sp<sub>1</sub>,Ge<sub>2</sub>,Sp<sub>2</sub>), ext:synonym(Ge<sub>2</sub>,Sp<sub>2</sub>,Ge,Sp).

$q_3(d_1,Lt,Ln,Ge,Sp) :-$  d<sub>2</sub>(Ge,Sp,Ct,Lt,Ln), Lt ≥ 33, Lt ≤ 42, Ln ≥ -124.3, Ln ≤ -115, ext:synonym('Manica','bradleyi',Ge<sub>1</sub>,Sp<sub>1</sub>), d<sub>4</sub>(Ge<sub>1</sub>,Sp<sub>1</sub>,Ge<sub>2</sub>,Sp<sub>2</sub>), ext:synonym(Ge<sub>2</sub>,Sp<sub>2</sub>,Ge,Sp).

$q_3(d_1,Lt,Ln,Ge,Sp) :-$  d<sub>3</sub>(Mp,Cf,Lt,Ln), Lt ≥ 33, Lt ≤ 42, Ln ≥ -124.3, Ln ≤ -115, ext:synonym('Manica','bradleyi',Ge<sub>1</sub>,Sp<sub>1</sub>), d<sub>4</sub>(Ge<sub>1</sub>,Sp<sub>1</sub>,Ge<sub>2</sub>,Sp<sub>2</sub>), ext:synonym(Ge<sub>2</sub>,Sp<sub>2</sub>, 'Manica', 'parasitica').

$q_3(d_1,Lt,Ln,Ge,Sp) :-$  d<sub>3</sub>(Mp,Cf,Lt,Ln), Lt ≥ 33, Lt ≤ 42, Ln ≥ -124.3, Ln ≤ -115, ext:synonym('Manica','bradleyi',Ge<sub>1</sub>,Sp<sub>1</sub>), d<sub>4</sub>(Ge<sub>1</sub>,Sp<sub>1</sub>,Ge<sub>2</sub>,Sp<sub>2</sub>), ext:synonym(Ge<sub>2</sub>,Sp<sub>2</sub>, 'Aphaenogaster', 'calderoni').

## 4 Summary and Future Work

The prototype described in this paper enables dataset-discovery queries and provides initial support for mixing external services with query-answering techniques. The prototype was written in Prolog and has an accompanying web interface for parameterizing (i.e., for selecting the geographic region, species, and parasitic relationship of interest) and displaying query results. The prototype also implements a simple description-logic reasoner for ontology classification, which is used in query answering (to find relevant mappings). As future work, we want to extend the prototype presented by (i) adding additional external services relevant to SEEK (e.g., for computing points within complex bounding boxes, incorporating gazetteers, adding additional synonym operations, etc.) and (ii) exploring techniques to enrich our framework for further exploiting arbitrary external services in query answering. As an example, consider the following semantic description for a dataset similar to d<sub>1</sub> and an externally defined service UTM2LatLon(Ux,Uy,Re,Zo,Lt,Ln) that converts UTM to latitude-longitude degree coordinates.

$d_5(Ge,Sp,Co,Ux,Uy,Re,Zo) :-$   
 Observation(O), value(O,Co), context(O,S), location(S,P), UTMPoint(P),  
 UTMx(P,Ux), UTMx(P,Uy), region(P,Re), zone(P,Zo), item(O,A), Abundance(A),  
 property(A,N), SciName(N), genus(N,Ge), species(N,Se).

In answering query  $q_1$ , we want to (1) return  $d_5$  as a relevant source, since UTM points can be converted to latitude-longitude points (using the external service), and (2) correctly insert the external service in the associated query rewriting. We are currently exploring *parameter dependency* specifications for this purpose, in which the domain and range of an external service are semantically described, e.g., using rules similar to the following.

**domain** UTM2LatLon(Ux,Uy,Re,Zo,Lt,Ln) :- UTMPoint(U), UTMx(Ux), UTMy(Uy),  
region(U,Re), zone(U,Zo).

**range** UTM2LatLon(Ux,Uy,Re,Zo,Lt,Ln) :- LatLonPoint(P), latDeg(P,Lt), lonDeg(P,Ln).

We believe that incorporating external services into data-integration architectures provides a powerful framework to support complex integration and transformation of scientific, and in particular, life-sciences data.

## References

1. W. Berensohn. The concept of “Potential Taxa” in databases. *Taxon*, vol. 44, 1995.
2. S. Bowers and B. Ludäscher. An ontology-driven framework for data transformation in scientific workflows. In *Proc. of Data Integration in the Life Sciences*, LNCS, vol. 2994, 2004.
3. B. Ludäscher, A. Gupta, and M. E. Martone. Model-based mediation with domain maps. In *Proc. of ICDE*, IEEE Computer Society, 2001.
4. B. Hölldobler and E. O. Wilson. *The Ants*. Harvard University Press, 1990.
5. N. F. Johnson. The Hymenoptera Name Server. [http://atbi.biosci.ohio-state.edu:210/hymenoptera/nomenclator.home\\_page](http://atbi.biosci.ohio-state.edu:210/hymenoptera/nomenclator.home_page)
6. M. Lenzerini. Data integration: A theoretical perspective. In *Proc. of PODS*, 2002.
7. A. Y. Levy, A. Rajaraman, and J. J. Ordille. Query-answering algorithms for information agents. In *Proc. of AAAI*, 1996.
8. D. R. B. Stockwell and D. P. Peters. The GARP modelling system: Problems and solutions to automated spatial prediction. *Intl. J. of Geographic Information Systems*, vol. 13, 1999.
9. A. Purvis and A. Hector. Getting the measure of biodiversity. *Nature*, vol. 405, 2000.
10. N. W. Paton, R. Stevens, P. Baker, C. A. Goble, S. Bechhofer, A. Brass. Query processing in the TAMBIS bioinformatics source integration system. In *Proc. of the SSDBM*, 1999.
11. T. Paterson and J. Kennedy. Approaches to storing and querying structural information in botanical specimen descriptions. To appear in *Proc. of BNCOD*, LNCS, July, 2004.