# Semantic linking of phenotypes and environments: A review

Anne E Thessen, Daniel E Bunker, Pier Luigi Buttigieg, Laurel D Cooper, Wasila M Dahdul, Sami Domisch, Nico M Franz, Pankaj Jaiswal, Carolyn J Lawrence-Dill, Peter E Midford, Christopher J Mungall, Martín J Ramírez, Chelsea D Specht, Lars Vogt, Rutger Aldo Vos, Ramona L Walls, Jeffrey W White, Guanyang Zhang, Andrew R Deans, Eva Huala, Suzanna E Lewis, Paula M Mabee

Understanding the interplay between environmental conditions and phenotypes is a fundamental goal of biology. Unfortunately, data that include observations on phenotype and environment are highly heterogeneous and thus difficult to find and integrate. One approach that is likely to improve the status quo involves the use of ontologies to standardize and link data about phenotypes and environments. Specifying and linking data in this manner will allow researchers to increase the scope and flexibility of large-scale analyses aided by modern computing methods. Investments in this area would advance diverse fields such as ecology, phylogenetics, and conservation biology. While several biological ontologies are well-developed, using them to link phenotypes and environments is rare because of gaps in ontological coverage and limits to interoperability among ontologies and disciplines. In this review, we present 1) use cases from diverse disciplines to illustrate questions that could be answered more efficiently using a robust linkage between phenotypes and environments, 2) two proof-of-concept analyses that show the value of linking phenotypes to environments in fishes and amphibians, and 3) two proposed example data models for linking phenotypes and environments using the extensible observation ontology (OBOE) and the Biological Collections Ontology (BCO) that can serve as a starting point for the development of a data model linking phenotypes and environments.

1  Semantic Linking of Phenotypes and Environments: A Review

2  Anne E. Thessen[1,2]*

3  Daniel E. Bunker[3]

4  Pier Luigi Buttigieg[4]

5  Laurel D. Cooper[5]

6  Wasila M. Dahdul[6]

7  Sami Domisch[7]

8  Nico M. Franz[8]

9  Pankaj Jaiswal[5]

10  Carolyn J. Lawrence-Dill[9]

11  Peter E. Midford[10]

12  Christopher J. Mungall[11]

13  Martín J. Ramírez[12]

14  Chelsea D. Specht[13]

15  Lars Vogt[14]

16  Rutger Aldo Vos[15]

17  Ramona L. Walls[16]

18  Jeffrey W. White[17]

19  Guanyang Zhang[8]

20  Andrew R. Deans[+18]

21  Eva Huala[+19]

22  Suzanna E. Lewis[+11]

23  Paula M. Mabee[+6]

24  [1]The Data Detektiv, Waltham, MA USA annethessen@gmail.com 443.225.9185

25  [2]The Ronin Institute for Independent Scholarship, Monclair, NJ USA

26  [3]Department of Biological Sciences, New Jersey Institute of Technology Newark, NJ USA

27   thedbunker@gmail.com

28  [4]Alfred-Wegener-Institut, Helmholtz-Zentrum für Polar- und Meeresforschung, Germany pbuttigi@mpi-

29   bremen.de

30  [5]Department of Botany and Plant Pathology, Oregon State University, Corvallis, OR USA

31   jaiswalp@science.oregonstate.edu and cooperl@science.oregonstate.edu

32  [6]Department of Biology, University of South Dakota, Vermillion, SD USA Wasila.Dahdul@usd.edu and

33    Paula.Mabee@usd.edu

34    [7]Department of Ecology and Evolutionary Biology, Yale University, New Haven, CT USA

35    sami.domisch@yale.edu

36    [8]Arizona State University, Tempe, AZ USA nico.franz@asu.edu and gzhang44@asu.edu

37    [9]Departments of Genetics, Development and Cell Biology and Agronomy, Iowa State University, Ames,

38    IA USA triffid@iastate.edu

39    [10]peter.midford@gmail.com

40    [11]Genomics Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA cjmungall@lbl.gov

41    and suzi@berkeleybop.org

42    [12]Division of Arachnology, Museo Argentino de Ciencias Naturales – CONICET, Buenos Aires,

43    Argentina ramirez@macn.gov.ar

44    [13]Departments of Plant and Microbial Biology & Integrative Biology, University and Jepson Herbaria,

45    University of California, Berkleley, Berkeley, CA, USA cdspecht@berkeley.edu

46    [14]Universität Bonn, Institut für Evolutionsbiologie und Ökologie, Bonn, Germany

47    lars.m.vogt@googlemail.com

48    [15]Naturalis Biodiversity Center, Leiden, The Netherlands rutger.vos@naturalis.nl

49    [16]iPlant Collaborative, University of Arizona, Tucson, AZ USA rwalls@iplantcollaborative.org

50    [17]USDA-ARS, US Arid Land Agricultural Research Center, Maricopa, AZ USA

51    Jeffrey.White@ARS.USDA.GOV

52    [18]Department of Entomology, Pennsylvania State University, University Park, PA, USA adeans@psu.edu

53    [19]Phoenix Bioinformatics, Redwood City, CA USA huala@phoenixbioinformatics.org

54

55    *corresponding author

56    [+]co-PIs of the Phenotype RCN grant

57

## Abstract

Understanding the interplay between environmental conditions and phenotypes is a fundamental goal of biology. Unfortunately, data that include observations on phenotype and environment are highly heterogeneous and thus difficult to find and integrate. One approach that is likely to improve the status quo involves the use of ontologies to standardize and link data about phenotypes and environments. Specifying and linking data in this manner will allow researchers to increase the scope and flexibility of large-scale analyses aided by modern computing methods. Investments in this area would advance diverse fields such as ecology,

66 phylogenetics, and conservation biology. While several biological ontologies are well-developed,

67 using them to link phenotypes and environments is rare because of gaps in ontological coverage

68 and limits to interoperability among ontologies and disciplines. In this review, we present 1) use

69 cases from diverse disciplines to illustrate questions that could be answered more efficiently

70 using a robust linkage between phenotypes and environments, 2) two proof-of-concept analyses

71 that show the value of linking phenotypes to environments in fishes and amphibians, and 3) two

72 proposed example data models for linking phenotypes and environments using the extensible

73 observation ontology (OBOE) and the Biological Collections Ontology (BCO) that can serve as

74 a starting point for the development of a data model linking phenotypes and environments.

75

76 **Introduction**

77 Phenotype is the expression of interactions between genotype and environment. This

78 relationship is fundamental to a wide range of biological research, from large-scale questions

79 about the effect of climate change on global ecosystems to small-scale questions involving

80 disease presentation in individual organisms. The urgency of such questions, coupled with the

81 "data deluge" (Hey, Tansley & Tolle, 2009), has motivated scientists to explore more efficient

82 ways to aggregate and explore life science data. Traditional methods of data dissemination,

83 publication, and deposition in stand-alone databases do not support the rapid, automated, and

84 integrative methods of data exploration needed to efficiently address pressing research priorities.

85 Two important barriers to understanding interactions between environment and

86 phenotype are the heterogeneity of terms and their imprecise definitions in data sets and

87 manuscripts. An ontology has the potential to tame this heterogeneity and allow researchers to

88 more efficiently query and manipulate, large-scale data sets (Fig. 1); however, several challenges

89 must be overcome before their benefits are realized. Historically, bio-ontologies first came into

90 popular use in the biomedical and model organism communities (Côté & Robboy, 1980;

91 Spackman, Campbell & Côté, 1997; Ashburner et al., 2000), but they are now being applied to

92 address much broader, comparative problem complexes (Mabee et al., 2007; Deans et al., 2015;

93 Dececchi et al., in press). A shift towards representing and reasoning over taxonomically diverse

94 phenotypes in an ontological framework has occurred in a period of less than 10 years and, not

95 surprisingly, brings about new semantic, computational, and even social challenges (Gerson,

96 2008). In this paper, we explore the difficulties of automated linking of environments and

97    phenotypes, review the current state-of-the-art, present use cases, and propose solutions to

98    frequently encountered problems.

99         Clearly representing the natural language descriptions of phenotypes and environments

100   with a set of ontologies is difficult, because natural language, while highly expressive, is often

101   semantically ambiguous and reliant on context (Sasaki and Putz, 2009; Seltmann et al., 2013).

102   Despite successes in developing standards within specific disciplines (e.g., Taylor et al., 2008),

103   standard vocabularies are rare and seldom widely or consistently used (Enke et al., 2012).

104   Individual scientists often have preferred terms with undocumented and highly nuanced

105   meanings (Chang & Schutze, 2006; see discussion in Huang et al., 2015). Further, there is a co-

106   evolution between natural language and ontologies, which can complicate the recording  of

107   provenance and backwards-compatibility (Seppälä, Smith & Ceusters, 2014; Ochs et al., 2015).

108   Thus, as it stands, a researcher wishing to perform a meta-analysis has to manually integrate data

109   sets, which often requires discussions with data providers to clarify meaning.

110        In addition to the intricacies of natural language used to describe phenotypes and

111   environments, the ontological representation of environments requires additional considerations.

112   Environments are often described using a collection of semantically complex (and often

113   ambiguous) terms that are applied differently across disciplines. Semantic representation of

114   terms such as "environment", "ecosystem", "habitat", "ecozone", "bioregion", and "biome" must

115   account for variable biological, ecological, geographic, geopolitical, and historical usage.  As a

116   result of this complexity, many specialized environmentally-themed terms such as

117   "Afrotropical", which are central to fields such as zoogeography and floristic science, are not yet

118   included in any ontology. Data about species interactions and behavior can be an essential

119   component in defining an organism's environment, but current ontological structures do not

120   include behavior regulation classes that can be tied to ecological processes (e.g., negative

121   regulation of foraging behavior by predation pressure). Environments can also be defined using a

122   more data-driven approach where a specific environment is defined as the intersection of

123   different factors, e.g., defining a desert as having a specific annual precipitation, temperature

124   range, and solar irradiation. In the field of plant science, some consider field management

125   practices (including, e.g., frequency irrigation or fertilizer application) to be a component of

126   environment, whereas others consider field management to differ from the environment because

127   these changes to conditions are not a part of the natural environment. Because of differences in

128  perspective, data are organized in different ways across multiple resources. These inconsistencies

129  in describing environments complicate analyses that might identify associations among

130  occurrences, phenotypes and environmental conditions.

131       In spite of these challenges, understanding relationships between the semantics of

132  environment and phenotype is fundamental for data integration and scientific progress in the

133  fields of conservation, agriculture, disease control, organismal development, and numerous

134  others in biology. Thus, there is a need for a more developed, flexible, and interlinked ontology

135  framework representing environments, phenotypes, and their interplay. This framework for

136  environments and phenotypes can allow automated inferencing over large, aggregated data sets,

137  as demonstrated for gene functions and biological processes (Ashburner et al., 2000). Below, we

138  present use cases that illustrate research questions that would benefit from semantically linking

139  environments and phenotypes and describe existing efforts working toward this goal.

140  **Background**

141       Within the field of informatics, classification strategies range from flat lists of terms, to

142  vocabularies, and ontologies. For example, a vocabulary might merely contain the terms "bone",

143  "leg", "femur" and their definitions. An ontology would further define these as classes and with

144  respect to their biological relationships by asserting that a "femur" is a type of "bone" and part of

145  the "leg". Moreover, such assertions are encoded in a standardized, machine-readable form. Thus

146  ontologies empower computers to reliably interpret and reason over these logical relationship

147  graphs. A well-known example of the technology's potential is provided by IBM's Watson

148  (Gliozzo et al., 2013). Ontologies are typically recorded in a syntax format such as the Web

149  Ontology Language (OWL; W3C OWL Working Group, 2012) or the Resource Description

150  Framework (RDF; http://www.w3.org/RDF/) that can be readily distributed and exchanged by

151  computers, thereby facilitating knowledge integration within a scientific community. For an

152  ontology to actually be useful to scientists, these same scientists must mutually agree upon,

153  develop, and nurture a shared collection of ontologies and the processes for maintaining it.

154       Over 40 ontologies and vocabularies have been created to describe environment and

155  phenotype (Table 1). Some of these resources are extended and refined through incorporation of

156  user and developer requests for new terms and cross-referencing terms to existing vocabularies.

157  Like software development, an essential aspect of ontology development is constant evaluation

158  through active use: describing data sets and asking key biological questions. To this end, a

159    number of groups are making the first inroads in the use of ontologies for studying the impact of

160    environment on phenotypes (e.g., The Encyclopedia of Life; Pafilis et al., 2015). The Minimum

161    Information for any Sequence (MIxS; Yilmaz et al., 2011) metadata checklist, a product of the

162    Genomic Standards Consortium (GSC; Field et al., 2011), does not specifically link phenotypes

163    to environments, but does include fields for describing environments using the Environment

164    Ontology (ENVO; www.environmentontology.org; Buttigieg et al., 2013) as well as the host

165    phenotype as part of the host-associated genome/metagenome environmental package. Although

166    MIxS recommends using terms from the Phenotypic Quality Ontology (PATO; Gkoutos et al.,

167    2004) in the host phenotype field, host phenotypes can be complex and could be described via a

168    mixture of phenotype ontologies (such as the Human Phenotype Ontology (Köhler et al., 2014)

169    or Mouse Phenotype Ontology (Gkoutos et al., 2004). The International Consortium for

170    Agricultural Systems Applications (ICASA) has built an infrastructure for combining genotype,

171    environment, and management data in agricultural analyses using a hierarchical data dictionary

172    (Hunt, White & Hoogenboom, 2001; White et al., 2013). This infrastructure is being integrated

173    in crop and climate modeling efforts, notably through the Agricultural Model Intercomparison

174    and Improvement Project (AgMIP), which promotes efforts to "simulate yield response to

175    climatic factors, abiotic factors, and genotypic variables" (http://research.agmip.org/). Oellrich et

176    al. (2015) recently developed a standardized method for describing and analyzing the phenotypes

177    associated with characterized mutant genes across species that includes environmental terms

178    from the Plant Environment Ontology (EO). Despite this progress, the available environment and

179    phenotype ontologies still contain major gaps in the coverage of their respective domains, and

180    significant investment is needed before data integration and analytics can be accomplished on a

181    large scale.

182

183    **Use Cases**

184         To communicate the importance of investing in environment and phenotype ontologies,

185    we present use cases drawn from several life science domains. These use cases represent the

186    types of research questions that either cannot currently be answered or can only be answered

187    with great difficulty.

188

189    **<u>Using Phenotype and Environment Ontologies in Ecology</u>**

190

191 **Coping with Climate Change in Conservation, Management, and Agriculture**

192 *Example Question:* Which species or crop varieties are projected to do well in my locality over

193 the next century?

194

195 *Background*:

196        Climate change is anticipated to affect environmental conditions with unprecedented

197 speed. Knowledge concerning the responses of ecological communities to these changes is very

198 limited: adaptation and migration are among numerous possibilities that must be considered.

199 Conservation and agricultural resources are also limited, so identifying and focusing

200 interventions on taxa that are less able to adapt can be very helpful. Besides commonly-used

201 projections based on species distributions models, another strategy for identifying at-risk species

202 is to assess their vulnerability based on their traits (i.e., phenotypes). By linking phenotypes to

203 specific environmental conditions, taxa or strains that are likely to thrive (or not) under those

204 conditions can be identified. For example, cataloguing phenotypes that are more prevalent

205 among organisms that live in hot and wet environments and detecting their presence in

206 organisms whose environments are warming and becoming more humid, allows some bearing on

207 the later organisms' ability to cope with such climate change. In agriculture, this can be used to

208 identify crop varieties that are likely to give higher or more stable yields under specific

209 conditions or wild relatives of crops that may possess useful traits. In conservation, similar

210 analyses can identify species at risk of extinction (Thormann et al., 2014). One system that hints

211 at performing this type of analysis currently is SemanticWildNET (Henderson, Khan & Hunter,

212 2007), which links data about birds and snakes to environmental conditions in Australia.

213

214 *Current Workflow*:

215        Steve works for a seed company that serves the southern Great Plains in the USA.

216 General Circulation Models (GCMs) project that over the next 25 to 30 years farmers in the

217 southern Great Plains will experience drier, warmer, and longer summers. His company wants to

218 start breeding sorghum hybrids that will perform well in these future conditions. Steve's

219 company has developed and phenotyped a wide range of parental lines that differ in yield

220 response and phenology under different environmental conditions, much of which is proprietary

221  data. He is able to find additional data sets that link geolocations and associated environmental

222  conditions to taxon phenotypes for crops (e.g., GRIN, the Germplasm Resources Information

223  Network; http://www.ars-grin.gov/npgs/) and link taxa to phenotypes (e.g. the TRY Plant Trait

224  Database, https://www.try-db.org/; Kattge et al., 2011). Environmental data sets that include

225  information about weather (NOAA), soil (USGS) and climate projections on a 1 km spatial grain

226  (WorldClim data set, Hijmans et al., 2005) are publicly available through government sources.

227      Steve decides that the best strategy for finding the top hybrid for a specific future habitat

228  is to manually link phenotypes to environmental conditions using the taxon name and location

229  (geographic coordinates) as a metadata bridge (Fig. 2). To work with the data, Steve must

230  download the files to a local machine. Because he does not have programming skills, he must

231  manually locate the specific data of interest from each data source and then make decisions about

232  appropriate integration using written documentation from the data provider. The data preparation

233  and integration takes six months of full time work.

234      When Steve finally has his data ready to analyze he must pick a statistical workflow and

235  software package that can identify phenotypes and environmental conditions that are observed

236  together. The next step would be to look at the climate projections to find the projected

237  environmental conditions his customers are likely to be facing and use these to identify the ideal

238  suite of phenotypes for that future climate regime. The final step would be to identify the taxa

239  that have the phenotypes in question.

240

241  *Future Workflow: Agriculture*

242      Steve works for a seed company that serves the USA. GCMs project that over the next 25

243  to 30 years farmers in the southern Great Plains will experience drier, warmer, and longer

244  summers. His company wants to start breeding sorghum hybrids that will perform well in these

245  future conditions. Steve's company has developed and phenotyped a wide range of parental lines

246  that differ in yield response and phenology, and these lines have been annotated using ontology

247  terms for traits (e.g., TO and PATO, Table 1) and the corresponding growth conditions (EO,

248  Table 1). Additionally, the habitat of each line (or its source organisms) is described by classes

249  from an environment ontology (ENVO, Table 1)

250      Steve queries his company's internal, semantically aware database for annotated records

251  corresponding to the lines his company has developed that have high yields under warmer

252 climatic conditions and when subject to drought. This gives him a list of candidate lines (i.e.,

253 phenotypes and genotypes) to use in development of new hybrids for the region. If the list of

254 these candidates does not provide sufficient resistance to high temperatures and drought, he may

255 choose to query a database containing information on the wild relatives of sorghum, along with

256 average rainfall and temperature data from the natural habitat of each species and/or annotations

257 describing their habitat using classes from an environment ontology. If necessary he will be able

258 to introgress genetic material encoding drought or high temperature tolerance from a wild

259 relative of sorghum into his breeding lines.

260

261 *Future Workflow: Wildlife Conservation*

262   Lupita is a park ranger that manages a coastal wildlife sanctuary. Some of the species in

263 her sanctuary are listed as threatened by the IUCN. According to the latest climate change

264 projections, her sanctuary is going to be hotter and wetter in 50 years. She has limited resources

265 to maintain the biodiversity in her sanctuary for the long term. After some thought, she decides

266 to identify at-risk species by comparing the traits of the organisms in her park with traits of

267 organisms that do not do well in hot, wet, coastal environments. Lupita searches a semantically

268 aware, publicly-available biology database and finds a list of traits for vertebrates whose habitats

269 do not include hot, wet, coastal environments and a list of traits for vertebrates with habitats that

270 do. Comparing differences between the two data sets gives a list of candidate traits which

271 suggest a taxon would be vulnerable to the projected climate regime. Searching for these traits

272 across the species in her sanctuary, Lupita identifies two species that are highly likely to fare

273 poorly in the projected climate, and she devotes resources to their conservation.

274

275 *Challenges Today*:

276   A large proportion of phenotype and environment data are part of the "long tail of dark

277 data" (Heidorn, 2008) that are not currently digital or discoverable. Although some phenotype,

278 environment, genotype, and climate projection datasets are available, they are difficult to find

279 and interrelate. These types of datasets can be cross-linked using space, time, or taxon, but the

280 formats of the different datasets can pose a challenge to integration (e.g., Reed, White & Brown

281 2003). In addition, metadata across disciplines, data types, and time periods are rarely consistent.

282 Key data items used for integration, such as taxonomic names, change over time and lead to

283  poorly linked data (Edwards et al., 2011; Giles, 2011; Page, 2008; Franz et al., 2015).

284

285  **<u>Using Phenotype and Environment Ontologies in Taxonomy</u>**

286  Connecting Specimen Phenotypes to Environment

287  *Example Question:* Which traits are common to beetles collected from deserts?

288

289  *Background*:

290       Natural history collections worldwide contain approximately two billion specimens

291  across various taxonomic groups (Ariño, 2010). Tens of millions of these specimens have their

292  phenomes at least partly described in the form of published taxonomic descriptions and may

293  have environmental data recorded on the specimen label or in a field notebook. Much of these

294  data have not been digitized and can be difficult to find and use. Connecting specimen-based

295  phenotype data to environmental information that describes where the specimen was isolated can

296  support predictive modeling of diversity and distribution.

297       The current massive digitization effort applied to collections is primarily done manually,

298  though efforts are being made to automate where possible (Barber, Lafferty & Landrum, 2013).

299  The environmental data associated with a specimen, typically a note on a specimen label, is

300  typically transcribed verbatim. If a curator wants to annotate a specimen with an environment or

301  habitat type or other metadata, the process of reading the information and relating it to an

302  ontology class is entirely manual. This is a very time-consuming workflow that involves

303  reconciling synonyms and disambiguating homonyms. Ideally, much of the manual labor of

304  reconciliation, disambiguation, and assignment would be shifted to a machine with curators

305  intervening only periodically.

306       A semantic model for representing specimen phenotypes has been developed (Balhoff,

307  Yoder & Deans, 2011) and applied to taxonomic descriptions (Mullins et al., 2012; Balhoff et al.,

308  2013). This model applies phenotype statements directly to specimens. Each specimen, residing

309  in an institutional collection, is associated with collecting event data, including where, when,

310  how, and by whom it was collected. The "where"-data could be connected to environment types

311  and other environmental data through semantic annotation using environment ontologies.

312

313  *Current Workflow*:

314     Kate wants to annotate insect specimens in a research collection with phenotypic and

315     environmental terms. All labels, published reports, and field notebooks concerning this collection

316     have been digitized and processed via optical character recognition ("OCRed"). She logs in as an

317     editor into the museum collections database that allows her to virtually access each specimen and

318     any associated documents. Kate begins working on the first specimen. The interface brings up

319     the label, an image of the specimen, the published description, and the relevant field notebook

320     page. An NLP-assisted algorithm within the interface reads the OCRed documents and highlights

321     environment-related terms in the text. She quickly reads the label and notebook near the

322     automated highlights, then searches ENVO to find the class that most accurately describes the

323     environmental conditions described by the collector. The environment listed in the notebook

324     does not match an existing ENVO class. She creates an issue in the ENVO issue tracker

325     (https://github.com/EnvironmentOntology/envo-p/issues/) requesting a new class that more

326     accurately describes the specimen environment. She does research on what the class should be

327     called and on the suggested definition. Kate will have to wait until someone at ENVO can

328     consider her request before completing the annotation.

329     After submitting the issue to ENVO, Kate reads the published description and goes back

330     and forth between PATO and the relevant insect anatomy ontology to find the classes she needs

331     to describe the specimen phenotypes and double-checks the classes by looking at the specimen

332     image (when possible). This process is very time consuming. When Kate is finished she adds the

333     relevant classes to the specimen database.

334     When Kate finally finishes her annotations, any user can query her museum website for

335     specimens that meet specific phenotypic and environmental constraints.

336

337     *Future Workflow*:

338     Kate wants to annotate insect specimens with phenotypic and environmental terms. All

339     materials concerning this collection have been digitized. Her museum has the cyberinfrastructure

340     that allows her to virtually access each specimen, bring up the related documents, and assign

341     relevant phenotypic and environmental terms through a point-and-click interface. Kate begins

342     working on the first specimen. The interface brings up the label, an image of the specimen, the

343     published description, and the relevant field notebook page. A text-mining tool highlights

344     relevant information in these sources and suggests classes from appropriate phenotype and

345   environment ontologies (e.g., ENVO). Kate agrees with the suggestions and clicks "approve".

346   For the next specimen, Kate agrees with the proposed phenotype classes, but does not see any

347   appropriate ENVO classes. She rejects the suggestions and is taken to a window that allows her

348   to browse ENVO for more appropriate classes. Still not satisfied, another click takes her to a

349   window that allows her to submit a request for a new class for which she suggests a definition

350   and relevant references. Kate is then taken back to the working environment and is presented

351   with the next specimen. Kate again agrees with the proposed phenotype classes, but the source

352   contains only a latitude and longitude for environment. The granularity offered by a lat/long

353   query does not capture microhabitats the insect may have been exposed to, but she decides that a

354   high-level description of the environment (by using biome or environmental feature classes in

355   ENVO) is preferable to providing no metadata. Kate opens a lat/long query window where she is

356   asked for a latitude, longitude, and date. Altitude is optional and depth is required for a lat/long

357   over water. This query returns environmental data relevant to the date that the specimen was

358   collected and the system suggests ENVO classes, some of which Kate agrees with. When Kate

359   finishes her annotations, any user can query her museum website for specimens that meet

360   specific phenotypic and environmental constraints.

361

362   *Challenges Today*:

363           Inconsistencies in geographic metadata associated with specimens are a major roadblock

364   in connecting phenotypes and environments (Vollmar, Macklin & Ford, 2010). Specimen

365   metadata are filled with ambiguous and synonymous terms with inconsistent granularity. For

366   example, the Plant Bug Inventory project database (http://research.amnh.org/pbi/; Schuh, 2012)

367   uses thousands of habitat names to describe the localities where insect species were collected,

368   including "cloud forest", "cloud forest with bamboo" and "cloud forest: oak trees, fern" (G.

369   Zhang pers. comm.). The documentation required to relate these terms to each other is currently

370   absent. In addition, high-level (but imprecise) locality information (e.g., "State College, Penn.")

371   is quite common for museum specimens and cannot be associated with fine-grained environment

372   types. Further, specimen labels often contain somewhat vague terms such as "neotropical" or

373   "mesohaline" that correspond to broad ecoregional definitions. According to Wikipedia,

374   mesohaline is defined as water that is between 5 and 18 salinity

375   (http://en.wikipedia.org/wiki/Salinity), but it is seldom clear whether a collector has intended a

376   precise definition such as this when writing the label. Thus, associating many specimens

377   currently in collections with well-defined environments may not be possible.

378   Some specimen metadata include a latitude and a longitude or a locality name, which

379   may be used to infer the environment, but environments change over time. For example, a

380   specimen may have been collected from a desert, which has since been paved over in the

381   expansion of a metropolitan area. Environments are also subject to cycles such as seasonal, diel,

382   or tidal. All of these factors make date and time important metadata. Annotating specimens in

383   more three-dimensional environments, such as the ocean or a mountain plateau, requires yet

384   another piece of information – depth or altitude.

385

386   **Using Phenotype and Environment Ontologies in Phylogeny**

387   **Reconstructing Ancestral Features and Habitats**

388   *Example Question:* Do species that have independently reduced or lost their eyes share common

389   environments now or in the past?

390

391   *Background*:

392   To infer the most probable features of a common ancestor given a phylogenetic tree and

393   the phenotypes of extant species, researchers utilize several well-developed parsimony and

394   likelihood methods. Similarly, the habitat preference of living species can be used to reconstruct

395   evolution of ecological niches. Connecting the phenotypic data from species with their habitat

396   and environmental data allows efficient analysis of these associations, allowing, for example, the

397   disentangling of evolutionary adaptation from other causes of phenotypic convergence.

398   Current methods of ancestral reconstruction rely on the uniform identification of a limited

399   number of environmental traits (e.g., habitats).  Users have parsimony, likelihood, and Bayesian

400   methods at their disposal for ancestral state reconstruction (e.g., Mesquite, Maddison &

401   Maddison, 2014; BEAST, BayesTraits, and R packages such as ape). These methods allow for

402   both discrete and continuous values.  For discrete characters, ancestral states are calculated from

403   the specific character states (e.g., environmental traits) found in the species.  For example, for a

404   clade of species that live in either "deep sea" or "underwater cavern" habitats, ancestral state

405   reconstruction is limited to these discrete habitats, i.e., the ancestor can be hypothesized to have

406   lived in either one or the other habitat.  However, an ontology can show that "deep sea" and

407    "underwater cavern" are both subtypes of an "aphotic marine environment", and thus this parent

408    term reveals this as a potential ancestral state for this clade.

409

410    *Current Workflow*:

411         Jane examines museum specimens of organisms belonging to a clade of freshwater fishes

412    which encompasses several hundred species. She discovers that the eyes vary in their level of

413    development: completely absent in some species, reduced in others, and fully developed in most.

414    After mapping this trait on a well-supported phylogeny, she concludes that eye reduction and

415    loss has occurred independently several times in this clade. This leads her to hypothesize that the

416    changes in eye development are associated with a species' habitat. She goes to the museum

417    databases and finds that the original descriptions of the collection sites for these specimens are

418    recorded as free text in the Darwin Core field "verbatimLocality"

419    (http://rs.tdwg.org/dwc/terms/verbatimLocality). She enters the "verbatimLocality" data into her

420    matrix of features mapped onto the phylogeny. Jane notices that several terms might be

421    synonymous and begins to research the specific definitions of the terms used and does her own

422    research into conditions at each locality. After one month of reconciling locality terms, she

423    begins to notice that species with reduced or absent eyes are all from subterranean environments.

424    She proceeds with her study, now examining other environmental factors or phenotypic traits

425    that might play a role in their shared habitat type.

426

427    *Future Workflow*:

428         While examining several hundred museum specimens of organisms belonging to a clade

429    of freshwater fishes, Jane discovers that the eyes vary in their level of development.  Mapping

430    this trait on a well-supported phylogeny shows that eye reduction and loss has occurred

431    independently several times in this clade. This leads her to hypothesize that the changes in eye

432    development are associated with a species' habitat. She goes to the museum databases and finds

433    that the original description of the place from where these specimens were collected was

434    recorded as free text in the Darwin Core field "verbatimLocality"

435    (http://rs.tdwg.org/dwc/terms/verbatimLocality), but the text is mapped to classes in an

436    environmental ontology such as ENVO.  She downloads these classes for all species and adds

437    them to the matrix of features that are mapped to the phylogeny.  She sees that species with

438 reduced or absent eyes are from localities variously described as "shallow pool in cave", "deep

439 water well", "deep phreatic habitat", and "swallow hole". A visualization tool allows her to see

440 the ontological classes which these descriptions have been mapped to as well as any shared

441 hierarchies or relations to other classes. She notices that these descriptions share "groundwater"

442 and an environmental material and "subterranean" as an environmental quality. She proceeds

443 with her study, now examining other environmental factors or phenotypic traits that might play a

444 role in their shared habitat type.

445

446 *Challenges Today*:

447 As in the other use cases, environmental ontologies must be provisioned to include the

448 classes relevant to a broad range of habitat types. Additionally, and similar to other use cases,

449 phenotypes of taxa that are represented in a computational format must be readily available. The

450 challenge unique to this use case is that methods of phylogenetic optimization that utilize

451 ontological relationships need to be developed. This will require consideration of the hierarchy

452 of class relationships such that the semantic similarity (Pesquita et al., 2009; Resnik, 1999)

453 among differing ancestral states at a particular node is taken into account when calculating the

454 appropriate assignment of a state to that node. Further, visualizations of the distribution of

455 phenotypic and environmental features on the tree that display, e.g., the most similar ontological

456 parent classes across nodes, need to be developed. An attempt to create an ancestral phenotype

457 ontology has previously been made by Ramírez & Michalik (2014).

458

459 **<u>Using Phenotype and Environment Ontologies in Behavioral Ecology</u>**

460

461 **Including Species Interactions in Habitat Assessments**

462 *Example Question:* How will this predatory wasp affect the spider population in my vegetable

463 garden?

464

465 *Background*:

466 Behavior is a phenotype that can be influenced by the presence or absence of other

467 organisms. The presence of other taxa can be just as important as abiotic features for determining

468 suitability of an environment for habitation by members of given species. An observation of a

469  taxon exhibiting a stress behavior is very different from an observation of the same taxon

470  exhibiting a feeding behavior. Changes in the ranges of organisms due to climate change or

471  accidental introduction is another way by which environments change and become more or less

472  suitable for specific phenotypes, such as feeding or courtship behaviors. These behaviors are

473  very important and when they are interrupted, can increase or decrease abundance of the affected

474  organism.

475  Current methods for retrieving behaviors that might be predictive of species interactions

476  mostly rely on published or unpublished ethograms and incidental comments in taxonomic

477  descriptions or experimental studies. There are databases of species interactions (Poelen, Simons

478  & Mungall, 2014), but these reflect interactions observed and reported in the literature, without

479  the behavioral content to make predictions about possible interactions resulting from the

480  introduction or range expansion of one or both species. Ideally, behavioral descriptions would

481  include specific environmental preferences as well as details of foraging, anti-predator, and

482  courtship behavior. The ability to make predictions of interactions would be an important

483  contribution when considering planned introductions or when setting priorities for preventing

484  unintentional spread.

485

486  *Current Workflow*:

487  Larry depends on his vegetable garden for food and on the spiders within it for pest

488  control. He frequently sees the jumping spider, *Phidippus clarus* Keyserling 1884, in the garden.

489  *P. clarus* is a widespread and common spider in the Eastern US (Edwards, 2004) and has been

490  demonstrated to be capable of controlling an experimental population of herbivorous insect pests

491  (Hoefler, Chen & Jakob, 2006). Larry hears from his local agricultural extension office that a

492  South American wasp that preys on spiders has been accidently introduced nearby. Should Larry

493  be concerned that the presence of the wasp will lead to more pests in his garden? Larry takes the

494  day off work to go to the local University library and asks a librarian to help him find

495  information about *P. clarus* and the South American wasp. Much of the information he needs is

496  in table format (ethogram) or in narrative text (comments in taxon descriptions and experimental

497  studies) and is difficult to decipher. The librarian makes him aware of a database of species

498  interactions that is easier to understand, but no data for *P. clarus* are available. At the end of the

499  day, Larry is still uncertain about the effect of the wasps on his garden spider population.

500

501 *Future Workflow*:

502        Larry depends on his vegetable garden for food and on the spiders within it, such as the

503 jumping spider *Phidippus clarus* Keyserling 1884, for pest control. Larry hears from his local

504 agricultural extension office that a South American wasp that preys on spiders has been

505 accidently introduced nearby. Should Larry be concerned that the presence of the wasp will lead

506 to more pests in his garden? Larry checks a gardening app on his mobile device that uses a

507 combination of ontologies and observation data to power a Q&A engine about nature in his area.

508 Through a simple user interface, he asks the app if the wasp is likely to affect the jumping spider

509 and whether there are additional potential consequences. Guided by the ontological structure

510 available in its back-end, the app states that 1) *P. clarus* is known to spend large amounts of time

511 on the tops and tips of plant shoots, and commonly lays its egg sacs near the tips of shoots

512 (Edwards 1980; Hill 2014), and 2) the wasp searches for prey on the tops and tips of plant

513 shoots. The inference engine used by the app are able to predict that the introduced wasp is likely

514 interact with Larry's population of *P. clarus* spiders. Because of where eggs are laid, this would

515 have the potential to interrupt *P. clarus* reproduction and thus reduce pressure on his garden

516 pests. With this information, Larry spends an hour making several wasp traps out of old plastic

517 bottles to place in his garden.

518

519 *Challenges Today*:

520        Environment ontologies currently do not explicitly incorporate species interactions in the

521 definition of their classes; however, an ontology for describing experimental conditions (EO)

522 does describe interactions between plants and other organisms in their environment. Many

523 environment ontologies, as they are currently structured, may not capture features relevant to

524 whether an environment will support specific behaviors, which can be very important data. Not

525 all taxa will engage in important behaviors in all environments, thus for many studies,

526 presence/absence data are not adequate. Creating a new set of ethological ontologies and

527 developing relations from their classes to those present in environmental ontologies has great

528 potential to address these issues, but requires significant effort to realize and maintain.

529

**DISCUSSION**

**Challenges**

532    The process of developing the Use Cases identified several major barriers to linking

533    phenotype and environment. These fall into two categories: challenges of coverage and

534    challenges of interoperability.

535

536    *Challenges of Coverage*

537    <u>Variable Granularity:</u>

538    Environmental data are reported with varying degrees of granularity that can take the

539    form of nested categories (e.g., continent – country – province – township – street), intervals

540    (e.g., ± 30 km), or significant digits (e.g., 5.236 vs 5.2). Some data sets, especially species

541    observations, include highly granular metadata specifying the exact location or exact conditions

542    under which a specimen was collected (such as collecting an insect from under tree bark or

543    collecting an organism in the presence of its predator). Although existing ontologies cover many

544    of the scales of interest, gaps prevent sufficient detail to capture all of the environmental data

545    provided in connection with collected specimens or published studies.  These are critical for

546    some taxa, such as insects collected from under bark (Jain & Balakrishnan, 2011). Currently,

547    such data are not discoverable due to the paucity of terms in existing ontologies and the lack of

548    easy-to-use tools that allow for semantic annotation with multiple ontologies.

549

550    <u>Terms and Definitions:</u>

551    One of the biggest challenges in creating ontologies for application to disciplines that

552    have a long history of published knowledge is the translation of the information in human-

553    readable narrative into a machine-readable form. Human language is very difficult for a machine

554    to understand largely because of its variability and nuance. Different terms (i.e., synonyms) can

555    be used to refer to the same concept, while a single term (i.e., homonyms) can refer to multiple,

556    different concepts. The human brain copes with this uncertainty by understanding context. One

557    way for a machine to cope with the variability of natural language is to provide it with an

558    ontology that includes synonymous terms; however, this can be difficult to maintain because

559    language evolves rapidly. Homonyms are a word-sense-disambiguation problem, which requires

560   heuristics about context to infer meaning; it is an active area of research (Zhan & Chen, 2011). A

561   homonym example applicable to environments is the term "scrubland", which means something

562   very different in California and South Africa. In this case, significant disambiguation could be

563   achieved by cross-referencing terms with geo-location or place names using resources such as

564   GAZ (Buttigieg et al., 2013), a gazetteer developed along ontological principles.

565

566   Incomplete Ontologies:

567   The development of ontologies in the biodiversity sciences has grown rapidly but is

568   relatively new, thus coverage is still small (Table 1). The OBO Foundry Library

569   (http://www.obofoundry.org/), a repository for biological ontologies, contains 22 ontologies

570   relevant to environments and phenotypes, with a total of 136,480 classes. Of these ontologies,

571   only one describes environments (ENVO) and one describes plant environmental conditions in

572   experimental treatments (EO). Eleven are phenotype or anatomy ontologies that cover specific

573   taxonomic groups, such as fungi (FYPO), animals (UBERON primarily for Chordates, with

574   other ontologies such as PORO for specific clades like Porifera (Thacker et al., 2013)), and

575   plants (TO) (Table 1).  Many other taxa, such as the microbial eukaryotes, do not have dedicated

576   ontologies. Furthermore, existing ontologies lack many key concepts required for application to

577   the many facets of biodiversity. This argues for the need of "living" ontologies (actively

578   maintained and highly-responsive to user requests) that can be updated continually and with

579   tools and services to allow users to request new classes and update existing classes with low

580   overhead. Ontology development is extremely time-consuming (Dahdul et al., 2015), and it must

581   be driven by scientific requirements, not by attempts to fully provision them *a priori*. Further,

582   provenance, i.e., the record of authorship involved in term development through persistent digital

583   identifiers such as ORCID (orcid.org), is a poorly developed feature in most ontologies, though

584   important for providing credit to contributors.

585

586   *Challenges of Interoperability*

587   Data Integration:

588   Linking environments, locations, and phenotypes will require interoperability between

589   several data types with the varying granularity used in biodiversity and geoscience. These

590   include data types from political and physical geography, coordinate systems, gazetteers, as well

591 as representations of environment and habitat. GeoNames has linked political geography and

592 some physical geography with coordinates (http://www.geonames.org/). A specimen with

593 coordinates can easily be linked to any number of political entities using the GeoNames API.

594 The same has not been accomplished for habitats; however, the components required to

595 accomplish this are falling into place. For example, the LifeMapper (Prajapati, 2009) and Map of

596 Life (Jetz, McPherson & Guralnick, 2012) projects use ecological niche modeling to map species

597 distributions based on environmental conditions. Additionally, the Encyclopedia of Life

598 TraitBank (http://eol.org/info/516) links taxa to their habitat type and phenotypic traits, but not to

599 geographic coordinates (Parr et al., in press). Once greater ontological representation of the link

600 between species and their environments is accomplished, robustly linking species' phenotypes to

601 their environments and locations become readily achievable.

602      In addition to spatial variation, environments show considerable variation over time and

603 often change over daily and seasonal cycles. This makes temporal data a key component for

604 meaningful integration. Environmental conditions measured at 14:00 can be very different from

605 those measured at 02:00 in the same location. The measurements made at the former, may not

606 apply to a specimen collected at the latter. In addition, an organism is rarely only exposed to

607 conditions measured at a single place and time. Some degree of integration is required to get a

608 complete picture of an environment associated with a phenotype (referred to as the "exposome"

609 in epidemiology, Wild, 2005).

610

611 <u>Ontology Legacy Alignment:</u>

612      The development of successful ontologies is often driven by a "bottom-up" community

613 approach. While this results in a product that is relevant for users, it can also result in multiple

614 partially overlapping ontologies, despite efforts to prevent duplication (e.g., Smith et al., 2007).

615 For improved integration and inferencing, overlapping ontologies need to be properly aligned

616 and those alignments need to be maintained over time. If not done properly, inferencing may be

617 inhibited or precluded altogether. This is a general problem that is not unique to environment or

618 phenotype ontologies (Cregan et al., 2005). A "top-down" approach to ontology development, in

619 which classes that constitute the top levels of a new ontology come from an existing domain or

620 upper-level ontology (e.g., CARO, UBERON, PO, BFO - Grenon & Smith, 2004; Haendel et al.,

621 2008; Mungall et al., 2012; Cooper et al., 2014), can result in a shared structure and

622 homogenized development across ontologies, although more specific classes will still require

623 alignment. Aligning ontologies manually is a large task and it is difficult to know the full

624 consequences of an alignment without testing (Ochs et al., 2015). The ability to support the

625 provenance of alignments and re-alignments can translate into trust and continued investment.

626 Numerous semi-automated tools for alignment have been developed (e.g., Granitzer et al., 2010;

627 Chen et al., 2014). Challenges include setting up proper relations between classes in different

628 ontologies such that the logical outcomes are valid and consistent (Franz & Peet 2009; Meilicke

629 & Stuckenschmidt 2009; Jiménez-Ruiz et al., 2009; Franz et al., 2015, N. Franz unpublished

630 data). The time and effort spent on maintaining alignments and interoperability can be eliminated

631 if shared community resources are instead developed (Dahdul et al., 2015). For example, several

632 independent anatomy ontologies for vertebrates [teleost (Dahdul et al., 2010); amphibian

633 (Maglia, Leopold & Pugener, 2007); vertebrate skeletal (Dahdul et al., 2012), and vertebrate

634 homologous organs (Niknejad et al., 2012)] were recently subsumed into UBERON, the

635 metazoan anatomy ontology (Haendel et al., 2014), and new content and development is now

636 focused on this single resource.

637

638 **Proof of Concept Demonstrations: Linking Environments and Phenotypes**

639 *Miniaturization in Fish*

640 Question: Has the evolution of miniaturization in fishes been driven by environmental variables?

641 Miniaturization is essentially the evolution of small body size and the associated set of

642 phenotypes, typically reduction or loss of structures. Authors have related this extreme change in

643 body size to organisms whose habitats include highly acidic waters, typical of peat bog or black

644 water habitats (Kottelat et al., 2006). As a proof of concept, we tested the hypothesis that

645 miniaturization is correlated with environmental variables. Using a list of miniaturized fishes and

646 their sister taxa extracted from the literature as input, we retrieved a phenotype X taxon synthetic

647 supermatrix from Phenoscape Knowledgebase (KB) (kb.phenoscape.org) using the Ontotrace

648 tool (Dececchi et al., in press). Using the common phenotype ontologies as a bridge, the KB

649 links evolutionary phenotypes of biodiverse taxa to candidate genes from model organisms.

650 Using the taxon names as input to GBIF, we created a list of 378 georeferenced observations

651 from museum specimen records (http://www.gbif.org/occurrence/download/0000659-

652 150211104439307; Fig. 3). These species' latitude and longitude occurrence records were

653     matched to the 1 km  HydroSHEDS hydrography (Lehner, Verdin & Jarvis, 2008) using a

654     horizontal distance tolerance of 3 km; they were then intersected with freshwater specific layers

655     by Domisch et al., (unpublished data). In this data set, the watershed of each 1 km stream reach

656     along the HydroSHEDS hydrography was delineated and then overlaid with climate (Hijmans et

657     al., 2005), topography (Lehner, Verdin & Jarvis, 2008), land cover (Tuanmu & Jetz, 2014), and

658     surface geology (USGS) layers. The differences in the habitat variables between miniatures and

659     non-miniatures were explored using a two-tailed t-test (Table 2). The results showed that

660     miniaturized fishes are found in warmer, wetter environments than their non-miniaturized

661     counterparts. New data layers are being developed to test specific phenotypic hypotheses related

662     to the habitats (e.g., pH, water flow) of miniaturized fishes (Domisch et al., unpublished data).

663     Specifically, the phenotypic data from the ontology-enabled matrix can be used to examine

664     correlations to environment with ontology-based miniaturized phenotypes (e.g., mandibular

665     sensory canal, absent; basibranchial 2 tooth, absent).

666

667     *Amphibian Reproduction*

668     <u>Question</u>: Which amphibians in my neighborhood are most likely to have their breeding

669     disrupted if a plan to drain a pond (the single source of year-round, standing freshwater) is

670     implemented?

671          The Encyclopedia of Life links environments associated with a given species' habitat and

672     phenotypes indirectly through taxon names. These data can be accessed and downloaded via

673     TraitBank (Quintero et al., 2014; Parr et al., in press). TraitBank uses Uniform Resource

674     Identifiers (URI), many from existing ontologies, as a controlled vocabulary for describing

675     characters and character states to facilitate large-scale data integration (Table 3). As proof of

676     concept, we queried TraitBank for breeding environment and developmental mode in 282

677     amphibian taxa. A Chi-Square Test was used to test for independence between habitat and

678     reproductive mode. The data suggested an important reproductive difference between

679     amphibians in aquatic and terrestrial habitats (Table 4). Ninety nine percent of the amphibians

680     with direct development breed in a terrestrial habitat. Ninety eight percent of the amphibians

681     with larval development (tadpoles) breed in an aquatic environment. This links the "larval

682     development" phenotype to the "freshwater" environment and the "direct development"

683     phenotype to the "terrestrial" environment. These data also suggest that a permanent freshwater

684   habitat is more important to amphibians with paedomorphic development than those with

685   indirect development.

686

687         These examples provide demonstrations of the value of linking phenotype to environment

688   and demonstrate how these links can be made with existing tools. More complicated research

689   questions are likely to require more nuanced linking for several reasons. First, phenotypes

690   frequently vary within a species; one cannot assume that every member of a species has the same

691   phenotype. In these two examples, we chose traits that were consistent across all members of a

692   species (miniaturization and developmental mode). In the miniaturization example, this allowed

693   addition of the GBIF query results to the Phenoscape Knowledgebase results. Second, an

694   organism's life style (ambush predator, nocturnal frugivore, etc.) within an environment is

695   deeply rooted to its phenotypic composition. For example, a visual predator in an environment

696   with low-light conditions may have a large eye phenotype while a scavenger in the same

697   environment may have a small eye phenotype. Trying to connect an eye size phenotype to this

698   environment would have to be clarified by including the ecological role of the taxon in a given

699   ecosystem. Third, scale can be important. Taxa of very different sizes can experience the same

700   environment in very different ways. For example, a soil protist will experience a forest

701   environment differently than a vascular plant. Despite these challenges, the highly simplified fish

702   and amphibian examples above still demonstrate the results of linking phenotypes and

703   environments with existing data and tools.

704

705   **Knowledge representation**

706         Despite the challenges of coverage and interoperability, we can demonstrate some basic

707   models linking phenotypes to environments using existing ontologies (Figs. 4,5). The Extensible

708   Observation Ontology (OBOE) provides a basic knowledge graph for linked measurements

709   (Madin et al., 2007; Madin et al., 2008). This ontology has been described in detail elsewhere

710   (Madin et al., 2007, Madin et al., 2008). Briefly, the fundamental OBOE model is built around

711   an "Observation" class which is an observation of an "entity" and has one or more

712   "measurements". Observations can also have a context of other observations. Phenotypes and

713   environments can be linked by representing an organism observation with a location observation

714   as its context (Fig. 4A). OBOE can model categorical and numerical measurements (Fig. 4B).

715  Thus, a geolocation, a data point, or a country code can be added to a location observation that

716  provides context for an organism observation. OBOE allows the use of literals as instances,

717  meaning a measurement can have as a value a string or a URI, which can be helpful when a

718  needed URI does not exist.

719      Although OBOE is well suited for describing observations, it was not originally built to

720  manage information about specimens or taxa. The Biological Collections Ontology (BCO)

721  (Walls et al., 2014a; Walls et al., 2014b; Deck et al., 2015) offers an alternative way to link data,

722  based on ontology design principles from the Ontology for Biomedical Investigations (OBI;

723  Brinkman et al., 2010), but adapted for biodiversity science. A key element of BCO is the

724  difference between a specimen collection process, which has a material entity (i.e., specimen) as

725  output and an observing process, which has data as output. Deck at al. (2015) describes how

726  information about locations (e.g., coordinates or environmental conditions) and taxonomy (e.g.,

727  the identification process or species name) can be linked to specimens. A similar approach can

728  be used to link phenotypic data to observations of organisms in their environment. At its most

729  basic, the BCO (via OBI) represents the observing process as a type of assay (an OBI class).

730  Rather than representing taxonomic information as an observation, BCO has a class for

731  taxonomic identification process, which, like assay, is a subclass of OBI:planned process (Fig.

732  5A). Fig. 5B shows how the same data from Figure 4B would be mapped to instances of BCO

733  classes.

734      OBOE and BCO were developed for different uses cases and therefore have different

735  approaches to representing observations. Nonetheless, there is significant overlap between the

736  two ontologies (e.g., OBOE's observation corresponds closely to BCO's observing process), and

737  ongoing efforts to align them are likely to lead to a harmonized model that can work for many

738  different use cases.

739

740  **Summary**

741      Providing data structures that improve integration of biological data is necessary for

742  efficiently addressing complex research questions. The link between phenotype and environment

743  is fundamental to research in taxonomy, ecology, and phylogenetics; its relevance extends to the

744  biomedical domain. One way to create this link is through the use of extensible ontologies

745  designed to work across different data types, such as OBOE or BCO in combination with ENVO

746 and other trait ontologies. Despite recent advances, significant challenges remain. We

747 recommend the following steps to increase interoperability between phenotype and environment

748 data:

749 Make it easy to contribute to existing ontologies.

750     The existing suite of ontologies is not adequate for linking phenotypes and environments

751 across the tree of life. To address this, new classes need to be added to extend and improve

752 existing taxonomy, phenotype, and envivonmental ontologies. Some ontologies have well-

753 developed pathways for submitting new classes and editing existing classes and resources to

754 respond to requests quickly (e.g., Gene Ontology), but frequently the social processes of

755 validating ontologies are not a part of the ontology platform.

756 Georeference environments with temporal considerations.

757     Many taxon observations are accompanied by geographical coordinates, collection date

758 and time, but lack adequate environmental descriptions. While services exist that can translate

759 coordinates into a municipality, retrieving environmental information using geographic

760 coordinates is not yet possible across the globe. In addition, because environments are dynamic,

761 temporal information should be used to filter results. A service is needed that can take

762 spatiotemporal information and return data concerning environmental conditions and ontology

763 classes corresponding to environment types. Map of Life can provide some data corresponding to

764 coordinates in some areas, but ontology classes are not yet available.

765 Organize research communities that share common resources.

766     Ontologies rely on community support, driven by scientific questions, to be relevant.

767 Communities of experts can be organized around workshops co-occuring at conferences and

768 funded through programs such as the National Science Foundation's Research Coordination

769 Network. Significant progress on discipline-specific ontologies has been made through the use of

770 targeted workshops (e.g., Yoder et al., 2010).

771

772 **Acknowledgments**

790

791

792

793

794

795 **Glossary**

796 **collecting event** – The process of specimen collection that occurs at a specific time and place.

797 **cyberinfrastructure** – The technological framework of interconnected databases and computers

798     across institutions that enable and support advanced, large-scale scientific research.

799 **Darwin Core** – A standard reference of terms related to biological diversity, in particular taxa

800     and their occurrences. Darwin Core was created to facilitate sharing of biodiversity

801     information.

802 **GCM (General Circulation Model)** – From Wikipedia: A general circulation model (GCM), a

803     type of climate model, is a mathematical model of the general circulation of a planetary

804     atmosphere or ocean and based on the Navier–Stokes equations on a rotating sphere with

805     thermodynamic terms for various energy sources (radiation, latent heat).

806 **genotype.** The genetic makeup or set of genes of an organism.

807 **georeferenced** – Observations or specimen collection records that are associated with locality

808     information (e.g., latitude and longitude).

809 **human readable** – Information that is presented in a format that can be understood by a human.

810 **inferencing** – Performed by software programs ("reasoners") that deduce logically consistent

811     statements implied by the entities and relations asserted in an ontology or database.

812 **knowledgebase** – A database of interconnected information.

813 **machine readable** – Information stored in a data format that can be understood by a computer.

814 **Machine Learning (ML)** – A type of artificial intelligence in which software programs have the

815     ability to learn (make decisions or data predictions) without being explicitly programmed

816     when given new data.

817 **meta-analysis** – A statistical analysis of data that is combined from independently conducted

818     research studies.

819 **NLP (Natural Language Processing)** – Methods used in computer programs to understand and

820     extract data from natural (human) language.

821 **ontology** – A set of  defined terms (classes, concepts) and the relations between them that

822     represent the knowledge of a particular domain. Terms in an ontology are related in a

823     directed, acyclical graph.

824 **OCR (Optical Character Recognition)** - automated conversion of images of text into machine-

825     readable text

826    **OWL (Web Ontology Language)** – The name encompassing the set of web-based languages

827      used for ontology building supported by the World Wide Web Consortium (WC3)

828      international standards body and based on the rules of formal semantics.

829    **phenome** – The entirety of an organism's phenotypic traits.

830    **phenotype** – One or more observable characteristics of an organism.

831    **provenance -** History of data and its place of origin.

832    **RDF (Resource Description Framework)** – A family of World Wide Web Consortium (WC3)

833      specifications originally designed as a metadata model and generally used to model

834      information in knowledge management applications

835    **semantic** – of or relating to meaning or context.

836    **semantic annotation** – The act of adding (i.e. 'tagging') information artifacts such as images,

837      free-text anatomical descriptions, or specimen collection records, with classes from an

838      ontology or similar resource which represents their meaning in a machine-readable fashion.

839    **specimen** – A whole organism or part of an organism preserved in a collection.

840    **taxonomic description** – Natural language description of a taxonomic group, typically includes

841      phenotypic characters such as morphology and behavior.

842    **URI (Uniform Resource Identifier) –** A string of characters used to identify a resource that

843    enables interactions with representations of the resource over the internet.

844    **vocabulary** – Flat list of terms that can be used to classify data. These terms are not explicitly

845    related to one another.

**References**

846

847 Allen V, Batello C, Berretta E, Hodgson J, Kothmann M, Li X, McIvor J, Milne J, Morris C,

848 Peeters A, Sanderson M, The Forage and Grazing Terminology Committee. 2011. An

849 international terminology for grazing lands and grazing animals. *Grass and Forage Science*

850 66:2–28. doi: 10.1111/j.1365-2494.2010.00780.x

851 Ariño A. 2010. Approaches to estimating the universe of natural history collections data.

852 *Biodiversity Informatics* 7:82-92. http://dx.doi.org/10.17161/bi.v7i2.3991

853 Arnaud E, Cooper L, Shrestha R, Menda N, Nelson RT, Matteis L, Skofic M, Bastow R, Jaiswal

854 P, Mueller L, McLaren G. 2012. Towards a reference plant trait ontology for modeling

855 knowledge of plant traits and phenotypes. In: KEOD 2012 – Proceedings of the

856 International Conference on Knowledge Engineering and Ontology Development, 220-225.

857 http://wrap.warwick.ac.uk/id/eprint/59831

858 Ashburner M, Ball C, Blake J, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight

859 SS, Eppig JT, Harris MA, Hill DP, Issel-Tarver L, Kasarskis A, Lewis S, Matese JC,

860 Richardson JE, Ringwald M, Rubin GM, Sherlock G. 2000. Gene Ontology: Tool for the

861 unification of biology. *Nature Genetics* 25:25-29. doi:10.1038/75556

862 Balhoff JP, Yoder M, Deans A. 2011. Linking semantic phenotypes to character matrices and

863 specimens. In *Biodiversity Information Standards (TDWG) annual meeting*. New Orleans.

864 Balhoff JP, Mikó I, Yoder M, Mullins PL, Deans AR. 2013. A semantic model for species

865 description, applied to the ensign wasps (Hymenoptera: Evaniidae) of New Caledonia.

866 *Systematic Biology* 62:639-659. doi: 10.1093/sysbio/syt028

867 Barber A, Lafferty D, Landrum L. 2013. The SALIX Method: A semi-automated workflow for

868 herbarium specimen digitization. *Taxon* 62:581-590. doi: http://dx.doi.org/10.12705/623.16

869 Brinkman R, Courtot M, Derom D, Fostel JM, He Y, Lord P, Malone J, Parkinson H, Peters B,

870 Rocca-Serra P, Ruttenberg A, Sansone SA, Soldatova LN, Stoeckert CJ Jr, Turner JA,

871 Zheng J, The OBI Consortium. 2010. Modeling biomedical experimental processes with

872 OBI. *J. Biomedical Semantics* 1:S7. doi: http://www.jbiomedsem.com/content/1/S1/S7

873 Buttigieg PL, Morrison N, Smith B, Mungall CJ, Lewis SE, ENVO Consortium. 2013. The

874 environment ontology: contextualising biological and biomedical entities. *Journal of*

875 *Biomedical Semantics* 4:43. doi:http://www.jbiomedsem.com/content/4/1/43

876 Ceusters W. 2012. An information artifact ontology perspective on data collections and

877      associated representational artifacts. *Studies in Health Technology and Informatics* 180:68-

878      72.

879 Chang JT, Schutze H. 2006. Abbreviations in biomedical text. In: *Text Mining for Biology and*

880      *Biomedicine* 99–119.

881 Chen M, Yu S, Franz N, Bowers S, Ludäscher B. 2014. Euler/X: A toolkit for logic-based

882      taxonomy integration. arXiv preprint arXiv:1402.1992

883 Cooper L, Elser JL, Preece J, Arnaud E, Stevenson DW, Todorovic S, Zhang E, Mungall C,

884      Smith B, Jaiswal P. 2014. Common reference ontologies for plant biology (cROP): A

885      platform for integrative plant genomics [abstract no. W819]. *Plant & Animal Genome XXII*

886      Available at: https://pag.confex.com/pag/xxii/webprogram/Paper9799.html

887 Côté RA, Robboy S. 1980. Progress in medical information management: Systematized

888      Nomenclature of Medicine (SNOMED). *Journal of the American Medical Association*

889      243:756-762. doi:10.1001/jama.1980.03300340032015

890 Cregan A, Mochol M, Vrandecic D, Bechhofer S. 2005. Pushing the limits of OWL, rules, and

891      Protégé. In: Grau B, Horrocks I, Parsia B, eds. *Proceedings of the OWLED\*05 Workshop*

892      *on OWL, Experiences and Directions*. CEUR-WS.

893 Dahdul WM, Balhoff JP, Engeman J, Grande T, Hilton EJ, Kothari CR, Lapp H, Lundberg JG,

894      Midford PE, Vision TJ, Westerfield M, Mabee PM. 2010. Evolutionary characters,

895      phenotypes and ontologies: Curating data from the systematic biology literature. *PloS ONE*

896      5:e10708. doi: 10.1371/journal.pone.0010708

897 Dahdul WM, Balhoff JP, Blackburn DC, Diehl AD, Haendel MA, Hall BK, Lapp H, Lundberg

898      JG, Mungall CJ, Ringwald M, Segerdell E, Van Slyke CE, Vickaryous MK, Westerfield M,

899      Mabee PM. 2012. A unified anatomy ontology of the vertebrate skeletal system. *PloS ONE*.

900      7:e51070. doi: 10.1371/journal.pone.0051070

901 Dahdul WM, Dececchi TA, Ibrahim N, Lapp H, Mabee PM. 2015. Moving the mountain:

902      analysis of the effort required to transform comparative anatomy into computable anatomy.

903      *Database*. 2015:bav040. doi:10.1093/database/bav040

904 Deans AR, Lewis SE, Huala E, Anzaldo SS, Ashburner M, Balhoff JP, Blackburn DC, Blake JA,

905      Burleigh JG, Chanet B, Cooper LD, Courtot M, Csösz S, Cui H, Dahdul W, Das S,

906      Dececchi TA, Dettai A, Diogo R, Druzinsky RE, Dumontier M, Franz NM, Friedrich F,

907    Gkoutos GV, Haendel M, Harmon LJ, Hayamizu TF, He Y, Hines HM, Ibrahim N, Jackson
908        LM, Jaiswal P, James-Zorn C, Köhler S, Lecointre G, Lapp H, Lawrence CJ, Le Novère N,
909        Lundberg JG, Macklin J, Mast AR, Midford PE, Mikó I, Mungall CJ, Oellrich A, Osumi-
910        Sutherland D, Parkinson H, Ramírez MJ, Richter S, Robinson RN, Ruttenberg A, Schulz
911        KS, Segerdell E, Seltmann KC, Sharkey MJ, Smith AD, Smith B, Specht CD, Squires RB,
912        Thacker RW, Thessen AE, Fernandez-Triana J, Vihinen M, Vize PD, Vogt L, Wall CE,
913        Walls RL, Westerfeld M, Wharton RA, Wirkner CS, Woolley JB, Yoder MJ, Zorn AM,
914        Mabee P. 2015. Finding our way through phenotypes. *PLoS Biology* 13:e1002033.
915        doi:10.1371/journal.pbio.1002033
916    Dececchi TA, Balhoff JP, Lapp H, Mabee PM. in press. Toward synthesizing our knowledge of
917        morphology: Using ontologies and machine reasoning to extract presence/absence
918        evolutionary phenotypes across studies. *Systematic Biology*. doi:10.1093/sysbio/syv031
919    Deck J, Guralnick R, Walls R, Blum S, Haendel M, Matsunaga A, Wieczorek J. 2015. Meeting
920        Report: Identifying practical applications of ontologies for biodiversity informatics.
921        *Standards in Genomic Sciences* 10:25. doi: 10.1186/s40793-015-0014-0
922    DiGiuseppe N, Pouchard L, Noy N. 2014. SWEET ontology coverage for earth system sciences.
923        *Earth Science Informatics* 7:249–264. doi: 10.1007/s12145-013-0143-1
924    Edwards G. 2004. Revision of the jumping spiders of the genus *Phidippus* (Araneae: Salticidae).
925        *Occasional papers of the Florida State Collection of Arthropods* 11:1–156.
926    Edwards G. 1980. Taxonomy, ethology, and ecology of *Phidippus* (Araneae: Salticidae) in
927        eastern North America. D. Phil. Thesis, University of Florida. Available at:
928        https://archive.org/details/taxonomyethology00edwa.
929    Edwards P, Mayernik M, Batcheller A, Bowker GC, Borgman CL. 2011. Science friction: Data,
930        metadata, and collaboration. *Social Studies of Science*. 41:667–690.
931        doi:0.1177/0306312711413314
932    Eilbeck K, Lewis SE, Mungall CJ, Yandell M, Stein L, Durbin R, Ashburner M. 2005. The
933        Sequence Ontology: A tool for the unification of genome annotations. Genome Biology
934        6:R44. doi:10.1186/gb-2005-6-5-r44
935    Enke N, Thessen AE, Bach K, Bendix J, Seeger B, Gemeinholzer B. 2012. The user's view on
936        biodiversity data sharing. *Ecological Informatics* 11:25-33.
937        http://dx.doi.org/10.1016/j.ecoinf.2012.03.004

938   Field D, Amaral-Zettler L, Cochrane G, Cole JR, Dawyndt P, Garrity G, Gilbert J, Glöckner FO,

939       Hirschman L, Karsch-Mizrachi I, Klenk HP, Knight R, Kottmann R, Kyrpides N, Meyer F,

940       Gil IS, Sansone SA, Schriml LM, Sterk P, Tatusova T, Ussery DW, White O, Wooley J.

941       2011. The Genomic Standards Consortium. *PLoS Biology* 9:e1001088.

942       doi:10.1371/journal.pbio.1001088

943   Franz N, Peet R. 2009. Towards a language for mapping relationships among taxonomic

944       concepts. *Systematics and Biodiversity* 7:5-20. doi: 10.1017/S147720000800282X

945   Franz NM, Chen M, Yu S, Kianmajd P, Bowers S, Ludäscher. 2015. Reasoning over taxonomic

946       change: Exploring alignments for the *Perelleschus* use case. *PLoS ONE* 10:e0118247.

947       doi:10.1371/journal.pone.0118247

948   Giglio M, Mungall C, Uetz P, Yin L, Goll J, Siegele D, Chibucos M, Hu J. 2009. Development

949       of an Ontology of Microbial Phenotypes (OMP). *Nature Precedings*.

950       doi:10.1038/npre.2009.3639.1

951   Giles JRA. 2011. Geoscience metadata - no pain, no gain. In: Sinha AK, Arctur D, Jackson eds.

952       *Societal Challenges and Geoinformatics*. Geological Society of America, 29–33.

953   Gkoutos G, Green E, Mallon A, Hancock JM, Davidson D. 2004. Using ontologies to describe

954       mouse phenotypes. *Genome Biology* 6:R8. doi:10.1186/gb-2004-6-1-r8

955   Gliozzo A, Biran O, Patwardhan S, McKeown K. 2013. Semantic Technologies in IBM Watson.

956       In: *Proceedings of the Fourth Workshop on Teaching NLP and CL*. Sofia, Bulgaria:

957       Association for Computational Linguistics, 85–92. Available at:

958       http://www.aclweb.org/anthology/W13-3413.

959   Granitzer M, Sabol V, Onn K, Lukose D, Tochtermann K. 2010. Ontology Alignment—A

960       Survey with Focus on Visually Supported Semi-Automatic Techniques. *Future Internet*

961       2:238–258. doi:10.3390/fi2030238

962   Grenon P, Smith B. 2004. SNAP and SPAN: Towards dynamic spatial ontology. *Spatial*

963       *Cognition & Computation: An Interdisciplinary Journal* 4:69-104.

964       doi:10.1207/s15427633scc0401_5

965   Haendel M, Neuhaus F, Osumi-Sutherland D, Mabee PM, Mejino JLV Jr, Mungall CJ, Smith B.

966       2008. CARO–the common anatomy reference ontology. In: *Anatomy Ontologies for*

967       *Bioinformatics*. London: Springer, 327–349.

968  Haendel M, Balhoff JP, Bastian F, Blackburn DC, Blake JA, Bradford Y, Comte A, Dahdul
969      WM, Dececchi TH, Druzinsky RE, Hayamizu TF, Ibrahim N, Lewis SE, Mabee PM,
970      Niknejad A, Robinson-Rechavi M, Sereno PC, Mungall CJ. 2014. Unification of multi-
971      species vertebrate anatomy ontologies for comparative biology in Uberon. *Journal of*
972      *Biomedical Semantics* 5:21. doi:10.1186/2041-1480-5-21
973  Harris MA, Lock A, Bähler J, Oliver SG, Wood V. 2013. FYPO: The fission yeast phenotype
974      ontology. Bioinformatics 29:1671-1678. doi:10.1093/bioinformatics/btt266
975  Hastings J, de Matos P, Dekker A, Ennis M, Harsha B, Kale N, Muthukrishnan V, Owen G,
976      Turner S, Williams M, Steinbeck C. 2013. The ChEBI reference and ontology for
977      biologically relevant chemistry: Enhancements for 2013. Nucleic Acids Research 41:D456-
978      D463. doi:10.1093/nar/gks1146
979  Heidorn PB. 2008. Shedding Light on the Dark Data in the Long Tail of Science. *Library Trends*
980      57:280–299.
981  Henderson M, Khan I, Hunter J. 2007. Semantic WildNET: An ontology based biogeographical
982      system, Available at:
983      http://www.itee.uq.edu.au/eresearch/filething/files/get/projects/ecoportalqld/SemWildNET.
984      pdf.
985  Hey T, Tansley S, Tolle K. 2009. *The fourth paradigm: data-intensive scientific discovery*,
986      Seattle: Microsoft Research.
987  Hijmans R, Cameron SE, Parra JL, Jones PG, Jarvis A. 2005. Very high resolution interpolated
988      climate surfaces for global land areas. *International Journal of Climatology*. 25:1965–1978.
989  Hill D. 2014. Notes on the jumping spider *Phidippus clarus* Keyserling 1885 (Araneae:
990      Salticidae: Dendryphantinae). *Peckhamia* 113.1:1–32.
991  Hoefler C, Chen A, Jakob E. 2006. The potential of a jumping spider, *Phidippus clarus*, as a
992      biocontrol agent. *Journal of Economic Entomology*. 99:432–436.
993      doi:http://dx.doi.org/10.1603/0022-0493-99.2.432
994  Huang F, Macklin J, Cui H, Cole HA, Endara L. 2015. OTO: Ontology Term Organizer. *BMC*
995      *Bioinformatics* 16:47. doi:10.1186/s12859-015-0488-1
996  Hunt LA, White JW, Hoogenboom G. 2001. Agronomic data: advances in documentation and
997      protocols for exchange and use. *Agricultural Systems* 70:477-492.

998   Jain M, Balakrishnan R. 2011. Microhabitat selection in an assemblage of crickets (Orthoptera:
999       Ensifera) of a tropical evergreen forest in Southern India. *Insect Conservation and Diversity*
1000      4:152–158.

1001  Jaiswal P, Ware D, Ni J, Chang K, Zhao W, Schmidt S, Pan X, Clark K, Teytelman L,
1002      Cartinhour S, Stein L, McCouch S. 2002. Gramene: Development and integration of trait
1003      and gene ontologies for rice. *Comparative and Functional Genomics* 3:132-136.
1004      doi:10.1002/cfg.156

1005  Jaiswal P, Avraham S, Ilic K, Kellogg EA, McCouch S, Pujar A, Reiser L, Rhee SY, Sachs MM,
1006      Schaeffer M, Stein L, Stevens P, Vincent L, Ware D, Zapata F. 2005. Plant Ontology (PO):
1007      A controlled vocabulary of plant structures and growth stages. Comparative Functional
1008      Genomics 6:388-397. doi:10.1002/cfg.496

1009  Jetz W, McPherson J, Guralnick R. 2012. Integrating biodiversity distribution knowledge:
1010      toward a global map of life. *Trends in Ecology & Evolution* 27:151–159.
1011      doi:10.1016/j.tree.2011.09.007

1012  Jiménez-Ruiz E, Grau B, Horrocks I, Berlanga R. 2009. Ontology Integration Using Mappings:
1013      Towards Getting the Right Logical Consequences. In: *The Semantic Web: Research and*
1014      *Applications*. Springer Berlin Heidelberg, 173–187.

1015  Kattge J, Diaz S, Lavorel S, Prentice C, Leadley P, Bonisch G, Garnier E, Westoby M, Reich
1016      PB, Wright IJ, Cornelissen JHC, Violle C, Harrison SP, van Bodegom PM, Reichstein M,
1017      Enquist BJ, Soudzilovskaia NA, Ackerly DD, Anand M, Atkin O, Bahn M, Baker TR,
1018      Baldocchi D, Bekker R, Blanco CC, Blonder B, Bond WJ, Bradstock R, Bunker DE,
1019      Casanoves F, Cavender-Bares J, Chambers JQ, Chapin FS, Chave J, Coomes D, Cornwell
1020      WK, Craine JM, Dobrin BH, Duarte L, Durka W, Elser J, Esser G, Estiarte M, Fagan WF,
1021      Fang J, Fernandez-Mendez F, Fidelis A, Finegan B, Flores O, Ford H, Frank D, Freschet
1022      GT, Fyllas NM, Gallagher RV, Green WA, Gutierrez AG, Hickler T, Higgins SI, Hodgson
1023      JG, Jalili A, Jansen S, Joly CA, Kerkhoff AJ, Kirkup D, Kitajima K, Kleyer M, Klotz S,
1024      Knops JMH, Kramer K, Kuhn I, Kurokawa H, Laughlin D, Lee TD, Leishman M, Lens F,
1025      Lenz T, Lewis SL, Lloyd J, Llusia J, Louault F, Ma S, Mahecha MD, Manning P, Massad
1026      T, Medlyn BE, Messier J, Moles AT, Muller SC, Nadrowski K, Naeem S, Niinemets U,
1027      Nollert S, Nuske A, Ogaya R, Oleksyn J, Onipchenko VG, Onoda Y, Ordonez J, Overbeck
1028      G, Ozinga WA, Patino S, Paula S, Pausas JG, Penuelas J, Phillips OL, Pillar V, Poorter H,

1029      Poorter L, Poschlod P, Prinzing A, Proulx R, Rammig A, Reinsch S, Reu B, Sack L,

1030      Salgado-Negre B, Sardans J, Shiodera S, Shipley B, Siefert A, Sosinski E, Soussana JF,

1031      Swaine E, Swenson N, Thompson K, Thornton P, Waldram M, Weiher E, White M, White

1032      S, Wright SJ, Yguel B, Zaehle S, Zanne AE, Wirth C. 2011. TRY - a global database of

1033      plant traits. *Global Change Biology* 17:2905-2935.

1034 Keyserling E. 1884-1885. Neue Spinnen aus Amerika. VI., Verhandlungen der k.k. zoologisch-

1035      botanischen Gesellschaft in Wien, Wien: 497. Available at:

1036      http://www.biodiversitylibrary.org/part/6448

1037 Kottelat M, Britz R, Hui T, Witte KE. 2006. *Paedocypris*, a new genus of Southeast Asian

1038      cyprinid fish with a remarkable sexual dimorphism, comprises the world's smallest

1039      vertebrate. *Proceedings of the Royal Society B* 273:895-899.

1040 Köhler S, Doelken SC, Mungall CJ, Bauer S, Firth HV, Bailleul-Forestier I, Black GC, Brown DL,

1041      Brudno M, Campbell J, FitzPatrick DR, Eppig JT, Jackson AP, Freson K, Girdea M, Helbig I, Hurst

1042      JA, Jähn J, Jackson LG, Kelly AM, Ledbetter DH, Mansour S, Martin CL, Moss C, Mumford A,

1043      Ouwehand WH, Park SM, Riggs ER, Scott RH, Sisodiya S, Van Vooren S, Wapner RJ, Wilkie AO,

1044      Wright CF, Vulto-van Silfout AT, de Leeuw N, de Vries BB, Washington NL, Smith CL,

1045      Westerfield M, Schofield P, Ruef BJ, Gkoutos GV, Haendel M, Smedley D, Lewis SE, Robinson

1046      PN. 2014. The Human Phenotype Ontology project: Linking molecular biology and disease through

1047      phenotype data. *Nucleic Acids Research* 42:D966-D974. doi:10.1093/nar/gkt1026

1048 Lehner B, Verdin K, Jarvis A. 2008. New global hydrography derived from spaceborne elevation

1049      data. *EOS, Transactions American Geophysical Union*. 89:93–94.

1050 Mabee PM, Ashburner M, Cronk Q, Gkoutos GV, Haendel M, Segerdell E, Mungall C,

1051      Westerfield M. 2007. Phenotype ontologies: the bridge between genomics and evolution.

1052      *Trends in Ecology and Evolution* 22:345-350. doi:10.1016/j.tree.2007.03.013

1053 Maddison W, Maddison D. 2014. Mesquite: a modular system for evolutionary analysis.

1054      Available at: http://mesquiteproject.org.

1055 Madin JS, Bowers S, Schildhauer MP, Krivov S, Pennington D, Villa F. 2007. An ontology for

1056      describing and synthesizing ecological observation data. *Ecological Informatics* 2:279–296.

1057      doi: 10.1016/j.ecoinf.2007.05.004

1058 Madin JS, Bowers S, Schildhauer MP, Jones MB. 2008. Advancing ecological research with

1059      ontologies. *Trends in Ecology & Evolution* 23:159-168.

1060 Maglia A, Leopold J, Pugener L. 2007. An anatomical ontology for amphibians. In: Altman R,
1061     Dunker A, Hunter L. eds. *Pacific Symposium on Biocomputing 2007*. World Scientific,
1062     367–378.

1063 Malone J, Holloway E, Adamsusiak T, Kapushesky M, Zheng J, Kolesnikov N, Zhukova A,
1064     Brazma A, Parkinson H. 2010. Modeling sample variables with an Experimental Factor
1065     Ontology. *Bioinformatics* 26:1112–1118. doi: 10.1093/bioinformatics/btq099

1066 Meilicke C, Stuckenschmidt H. 2009. An Efficient Method for Computing Alignment
1067     Diagnoses. In: *Web Reasoning and Rule Systems*. Springer Berlin Heidelberg, 182–196.

1068 Mullins P, Kawada R, Balhoff JP, Deans AR. 2012. A revision of *Evaniscus* (Hymenoptera,
1069     Evaniidae) using ontology-based semantic phenotype annotation. *ZooKeys* 223:1–38.
1070     doi:10.3897/zookeys.223.3572

1071 Mungall C, Torniai C, Gkoutos G, Lewis SE, Haendel MA. 2012. Uberon, an integrative multi-
1072     species anatomy ontology. *Genome Biology* 13:R5. doi: 10.1186/gb-2012-13-1-r5

1073 Niknejad A, Comte A, Parmentier G, Roux J, Bastian FB, Robinson-Rechavi M. 2012. vHOG, a
1074     multispecies vertebrate ontology of homologous organs groups. *Bioinformatics*. 28:1017–
1075     1020. doi: 10.1093/bioinformatics/bts048

1076 Ochs C, Perl Y, Geller J, Haendel MA, Brush M, Arabandi S, Tu S. 2015. Summarizing and
1077     visualizing structural changes during the evolution of biomedical ontologies using a Diff
1078     Abstraction Network. *Journal of Biomedical Informatics* 56:127-144.
1079     doi:10.1016/j.jbi.2015.05.018

1080 Oellrich A, Walls RL, Cannon EK, Cannon SB, Cooper L, Gardiner J, Gkoutos GV, Harper L,
1081     He M, Hoehndorf R, Jaiswal P, Kalberer SR, Lloyd JP, Meinke D, Menda N, Moore L,
1082     Nelson RT, Pujar A, Lawrence CJ, Huala E. 2015. An ontology approach to comparative
1083     phenomics in plants. *Plant Methods* 11:10. doi: 10.1186/s13007-015-0053-y

1084 Pafilis E, Frankild SP, Schnetzer J, Fanini L, Faulwetter S, Pavloudi C, Vasileiadou K, Leary P,
1085     Hammock J, Schulz K, Parr CS, Arvanitidis C, Jensen LJ. 2015. ENVIRONMENTS and
1086     EOL: Identification of Environment Ontology terms in text and the annotation of the
1087     Encyclopedia of Life. *Bioinformatics* 31:1872-1874. doi:10.1093/bioinformatics/btv045

1088 Page RDM. 2008. Biodiversity informatics: The challenge of linking data and the role of shared
1089     identifiers. *Briefings in Bioinformatics* 9:345–354.

1090    Parr C, Wilson N, Schulz K, Leary P, Hammock J, Rice J, Corrigan RJ Jr. in press. TraitBank:
1091        Practical semantics for organism attribute data. *Semantic Web*. Available at:
1092        http://www.semantic-web-journal.net/system/files/swj650.pdf.

1093    Pesquita C, Faria D, Falcao A, Lord P, Couto FM. 2009. Semantic similarity in biomedical
1094        ontologies. *PLoS Computational Biology* 5:e1000443.

1095    Poelen JH, Simons JD, Mungall CJ. 2014. Global Biotic Interactions: An open infrastructure to
1096        share and analyze species-interaction datasets. *Ecological Informatics* 24:148–159.
1097        doi:10.1016/j.ecoinf.2014.08.005

1098    Prajapati V. 2009. LIFEMAPPER: Mapping and Predicting the Distribution of Life with
1099        Distributed Computation: The Future of Biodiversity. *Archives of Applied Science Research*
1100        1:306–312.

1101    Quintero E, Thessen AE, Arias-Caballero P, Ayala-Orozco B. 2014. A statistical assessment of
1102        population trends for data deficient Mexican amphibians. *PeerJ* 2:e703.
1103        doi:https://dx.doi.org/10.7717/peerj.703

1104    Ramírez MJ, Michalik P. 2014. Calculating structural complexity in phylogenies using ancestral
1105        ontologies. *Cladistics* 30:635-649. doi:10.1111/da.12075

1106    Reed B, White M, Brown J. 2003. Remote sensing phenology. In: Schwartz M, ed. *Phenology:*
1107        *An Integrative Environmental Science*. Springer Netherlands, 365–381.

1108    Resnik P. 1999. Semantic similarity in a taxonomy: An information-based measure and its
1109        application to problems of ambiguity in natural language. *Journal of Artificial Intelligence*
1110        *Research* 11:95–130.

1111    Sasaki N, Putz FE. 2009. Critical need for new definitions of "forest" and "forest degradation" in
1112        global climate change agreements. *Conservation Letters* 2:226-232. doi:10.1111/j.1755-
1113        263X.2009.00067.x

1114    Schindelman G, Fernandes JS, Bastiani CA, Yook K, Sternberg PW. 2011. Worm Phenotype
1115        Ontology: Integrating phenotype data within and beyond the *C. elegans* community. B*MC*
1116        *Bioinformatics* 12:32. doi:10.1186/1471-2105-12-32

1117    Schuh R. 2012. Integrating specimen databases and revisionary systematics. *ZooKeys* 209:255–
1118        267.

1119  Seltmann K, Pénzes Z, Yoder M, Bertone MA, Deans AR. 2013. Utilizing descriptive statements
1120       from the Biodiversity Heritage Library to expand the Hymenoptera Anatomy Ontology.
1121       *PLoS ONE* 8:e55674. doi: 10.1371/journal.pone.0055674
1122  Seppälä S, Smith B, Ceusters W. 2014. Applying the realism-based ontology versioning method
1123       for tracking changes in the Basic Formal Ontology. In: Garbacz P, Kutz O, eds. *Formal*
1124       *Ontology in Information Systems*. IOS Press, 227-240.
1125  Shrestha R, Arnaud E, Mauleon R, Davenport GF, Hancock D, Morrision N, Bruskiewich R,
1126       McLaren G. 2010. Multifunctional crop trait ontology for breeders' data: Field book,
1127       annotation, data discovery and semantic enrichment of the literature. *AoB Plants*
1128       2010:plq008. doi: 10.1093/aobpla/plq008
1129  Smith CL, Eppig JT. 2009. The Mammalian Phenotype Ontology: Enabling robust annotation
1130       and comparative analysis. Wiley Interdisciplinary Reviews: Systems Biology and Medicine
1131       1:390-399. doi:10.1002/wsbm.44
1132  Smith B, Ashburner M, Rosse C, Bard J, Bug W, Ceusters W, Goldberg LJ, Eilbeck K, Ireland
1133       A, Mungall C, The OBI Consortium, Leontis N, Rocca-Serra P, Ruttenber A, Sansone SA,
1134       Scheuermann RH, Shah N, Whetzel PL, Lewis S. 2007. The OBO Foundry: Coordinated
1135       evolution of ontologies to support biomedical data integration. *Nature Biotechnology*
1136       25:1251-1255. doi:10.1038/nbt1346
1137  Spackman KA, Campbell KE, Côté RA. 1997. SNOMED RT: A reference terminology for
1138       health care. *Proceedings of the AMIA Annual Fall Symposium* 1997:640-644.
1139  Takhtajan A. 1986. *Floristic Regions of the World*, University of California Press.
1140  Taylor CF, Field D, Sansone SA, Aerts J, Apweiler R, Ashburner M, Ball CA, Binz PA, Bogue
1141       M, Booth T, Brazma A, Brinkman RR, Clark AM, Deutsch EW, Fiehn O, Fostel J, Ghazal
1142       P, Gibson F, Gray T, Grimes G, Hancock JM, Hardy NW, Hermjakob H, Julian RK, Kane
1143       M, Kettner C, Kinsinger C, Kolker E, Kuiper M, Le Novere N, Leebens-Mack J, Lewis SE,
1144       Lord P, Mallon AM, Marthandan N, Masuya H, McNally R, Mehrle A, Morrison N,
1145       Orchard S, Quackenbush J, Reecy JM, Robertson DG, Rocca-Serra P, Rodriguez H,
1146       Rosenfelder H, Santoyo-Lopez J, Scheuermann RH, Schober D, Smith B, Snape J,
1147       Stoeckert CJ, Tipton K, Sterk P, Untergasser A, Vandesompele J, Wiemann S. 2008.
1148       Promoting coherent minimum reporting guidelines for biological and biomedical

1149      investigations: the MIBBI project. *Nature Biotechnology* 26:889–896. doi:10.1038/nbt0808-

1150      889

1151 Thacker RW, Díaz MC, Kerner A, Vignes-Lebbe R, Segerdell E, Haendel MA, Mungall CJ.

1152      2013. The Porifera Ontology (PORO): enhancing sponge systematics with an anatomy

1153      ontology. *Journal of Biomedical Systematics* 5:39. doi:10.1186/2041-1480-5-39

1154 Thormann I, Parra-Quijano M, Endresen D, Rubio-Teso ML, Iriondo MJ, Maxted N. 2014.

1155      Predictive characterization of crop wild relatives and landraces. Biodiversity International.

1156      Available at http://www.bioversityinternational.org/e-library/publications/

1157 Tuanmu M, Jetz W. 2014. A global 1km consensus land cover product for biodiversity and

1158      ecosystem modelling. *Global Ecology and Biogeography* 23:1031–1045.

1159 USGS International Surface Geology

1160      http://certmapper.cr.usgs.gov/data/envision/index.html?widgets=geologymaps

1161 Vos RA, Biserkov J, Balech B, Beard N, Blissett M, Brenninkmeijer C, van Dooren T, Eades D,

1162      Gosline G, Groom QJ, Hamann TD, Hettling H, Hoehndorf R, Helleman A, Hovenkamp P,

1163      Kelbert P, King D, Kirkup D, Lammers Y, DeMeulemeester T, Mietchen D, Miller JA,

1164      Mounce R, Nicolson N, Page R, Pawlik A, Pereira S, Penev L, Richards K, Sautter G,

1165      Shorthouse DP, Tahtinen M, Weiland C, Williams AR, Sierra S. 2014. Enriched

1166      biodiversity data as a resource and service. *Biodiversity Data Journal* 2:e1125.

1167      doi:10.3897/BDJ.2.e1125

1168 Vollmar A, Macklin JA, Ford L. 2010. Natural history specimen digitization: challenges and

1169      concerns. *Biodiversity Informatics* 7:93-112. doi: http://dx.doi.org/10.17161/bi.v7i2.3992

1170 W3C OWL Working Group, 2012. *OWL 2 Web Ontology Language*, Available at:

1171      http://www.w3.org/TR/2012/REC-owl2-overview-20121211.

1172 Walls R, Deck J, Guralnick R, Baskauf S, Beaman R, Blum S, Bowers S, Buttigieg PL, Davies

1173      N, Endresen D, Gandolfo MA, Hanner R, Janning A, Krishtalka L, Matsunaga A, Midford

1174      P, Morrison N, Tuama EO, Schildhauer M, Smith B, Stucky BJ, Thomer A, Wieczorek J,

1175      Whitacre J, Wooley J. 2014a. Semantics in support of biodiversity knowledge discovery:

1176      An introduction to the biological collections ontology and related ontologies. *PLoS ONE*

1177      9:e89606.

1178 Walls R, Guralnick R, Deck J, Buntzman A, Buttigieg PL, Davies N, Denslow MW, Gallery RE,

1179      Parnell JJ, Osumi-Sutherland D, Robbins RJ, Rocca-Serra P, Wieczorek J, Zheng J. 2014b.

1180        Meeting report: Advancing practical applications of biodiversity ontologies. *Standards in*

1181        *Genomic Sciences* 9:17. doi:10.1186/1944-3277-9-17

1182    White J, Hunt L, Boote K, Jones JW, Koo J, Kim S, Porter CH, Wilkens PW, Hoogenboom G.

1183        2013. Integrated description of agricultural field experiments and production: The ICASA

1184        Version 2.0 data standards. *Computers and Electronics in Agriculture* 96:1–12.

1185        doi:10.1016/j.compag.2013.04.003

1186    Wild C. 2005. Complementing the genome with an "exposome": The outstanding challenge of

1187        environmental exposure measurement in molecular epidemiology. *Cancer Epidemiology,*

1188        *Biomarkers & Prevention* 14:1847. doi: 10.1158/1055-9965.EPI-05-0456

1189    Yilmaz P, Kottmann R, Field D, Knight R, Cole JR, Amaral-Zettler L, Gilbert JA, Karsch-

1190        Mizrachi I, Johnston A, Cochrane G, Vaughan R, Hunter C, Park J, Morrison N, Rocca-

1191        Serra P, Sterk P, Arumugam M, Bailey M, Baumgartner L, Birren BW, Blaser MJ, Bonazzi

1192        V, Booth T, Bork P, Bushman FD. 2011. Minimum information about a marker gene

1193        sequence (MIMARKS) and minimum information about any (x) sequence (MIxS)

1194        specifications. *Nature Biotechnology* 29:415-420. doi:10.1038/nbt.1823

1195    Yoder MJ, Mikó I, Seltmann KC, Bertone MA, Deans AR. 2010. A Gross Anatomy Ontology

1196        for Hymenoptera. *PLoS ONE* 5:e15991. doi: 10.1371/journal.pone.0015991

1197    Zhan J, Chen Y. 2011. Research on Word Sense Disambiguation. *Advanced Materials Research*

1198        181-182:337–342.

1199

1200

1201　Table 1: List of resources (vocabularies and ontologies) relevant to annotating phenotypes and

1202　environments.

| Name | Abbreviation | URL | Reference |
|---|---|---|---|
| AGROVOC | | 1 | |
| Behavioral Ontology | NBO | 2 | |
| Chemical Entities of Biological Interest | ChEBI | 3 | Hastings et al., 2013 |
| CMECS Habitat Classification | | 4 | |
| Crop Ontology | CO | 5 | Shrestha et al., 2010 |
| Eagle-i Resource Ontology | ERO | 6 | |
| EcoLexicon | | 7 | |
| Ecological Classifications NatureServe | | 8 | |
| Environment Ontology | ENVO | 9 | Buttigieg et al., 2013 |
| EUNIS Habitat Classification | | 10 | |
| Experimental Factor Ontology | EFO | 11 | Malone et al., 2010 |
| Exposure ontology | EXO | 12 | |
| Fission Yeast Phenotype Ontology | FYPO | 13 | Harris et al., 2013 |
| Flora Phenotype Ontology | FLOPO | 14 | Vos et al., 2014 |
| Floristic Regions of the World | | | Takhtajan, 1986 |
| Fungal gross anatomy | FAO | 15 | |
| Gazetteer | GAZ | 16 | |
| Gene Ontology | GO | 17 | Ashburner et al., 2000 |
| GeoNames | | 18 | |
| Getty Thesaurus of Geographic Names | | 19 | |
| Global Administrative Areas | GADM | 20 | |
| Human Phenotype Ontology | HP | 21 | Köhler et al., 2014 |
| Information Artifact Ontology | IAO | 22 | Ceusters, 2012 |
| International Consortium for Agricultural Systems Applications standards | ICASA | | White et al., 2013 |
| IUCN Habitats Classification Scheme | | 23 | |
| Mammalian phenotype | MP | 24 | Smith and Eppig, 2009 |
| Mapping European Seabed Habitats | MESH | 25 | |
| NASA GCMD keyword list for locations | | 26 | |
| Ontology of Biological Attributes | OBA | 27 | |
| Ontology of Biomedical Investigation | OBI | 28 | Brinkman et al., 2010 |
| Ontology of Microbial Phenotypes | OMP | 29 | Giglio et al., 2009 |
| Phenotype Quality Ontology | PATO | 30 | Gkoutos et al., 2004 |
| Plant Environment Ontology | EO | 31 | |
| Plant Ontology | PO | 32 | Jaiswal et al., 2005 |
| Plant Trait Ontology | TO | 33 | Jaiswal et al., 2002 Arnaud et al., 2012 |
| Population and Community Ontology | PCO | 34 | |
| Relation Ontology | RO | 35 | |
| Semantic Web for Earth and Environmental Terminology | SWEET | 36 | DiGiuseppe et al., 2014 |
| Sequence Ontology | SO | 37 | Elbeck et al., 2005 |

| | | | |
|---|---|---|---|
| Terminology of Grazing Lands and Grazing Animals | | | Allen et al., 2011 |
| Uber Anatomy Ontology | UBERON | 38 | Mungall et al., 2012; Haendel et al., 2014 |
| Worm Phenotype | WBPhenotype | 39 | Schindelman et al., 2011 |
| WWF Ecozones | | 40 | |

1203    1 http://aims.fao.org/agrovoc#.VG4QG_nF_ng
1204    2 https://code.google.com/p/behavior-ontology/
1205    3 https://www.ebi.ac.uk/chebi/
1206    4 https://marinemetadata.org/references/cmecshabitat
1207    5 http://pantheon.generationcp.org/index.php?option=com_content&task=section&id=7&Itemid=35
1208    6 https://www.eagle-i.net/
1209    7 http://ecolexicon.ugr.es/en/index.htm
1210    8 http://explorer.natureserve.org/classeco.htm
1211    9 http://www.environmentontology.org
1212    10 https://marinemetadata.org/references/eunishabitat
1213    11 http://www.ebi.ac.uk/efo/
1214    12 http://www.obofoundry.org/cgi-bin/detail.cgi?id=exo
1215    13 http://www.pombase.org/
1216    14 http://wiki.pro-ibiosphere.eu/wiki/Traits_Task_Group
1217    15 http://www.yeastgenome.org/fungi/fungal_anatomy_ontology/
1218    16 http://bioportal.bioontology.org/ontologies/GAZ
1219    17 http://geneontology.org/
1220    18 http://www.geonames.org/
1221    19 http://www.getty.edu/research/tools/vocabularies/tgn/index.html
1222    20 http://www.gadm.org/
1223    21 http://www.human-phenotype-ontology.org/
1224    22 https://code.google.com/p/information-artifact-ontology/
1225    23 http://www.iucnredlist.org/technical-documents/classification-schemes/habitats-classification-scheme-ver3
1226    24 http://www.informatics.jax.org/searches/MP_form.shtml
1227    25 http://www.emodnet-seabedhabitats.eu/
1228    26 https://marinemetadata.org/references/cfregions
1229    27 http://wiki.geneontology.org/index.php/Extensions/x-attribute
1230    28 http://obi-ontology.org/page/Main_Page
1231    29 http://microbialphenotypes.org/wiki/index.php/Main_Page
1232    30 http://obofoundry.org/wiki/index.php/PATO:Main_Page
1233    31 http://planteome.org/amigo/cgi-bin/crop_amigo/browse.cgi?
1234    32 http://www.plantontology.org/
1235    33 http://planteome.org/amigo/cgi-bin/crop_amigo/browse.cgi?
1236    34 https://github.com/PopulationAndCommunityOntology/pco
1237    35 https://github.com/oborel/obo-relations
1238    36 https://sweet.jpl.nasa.gov/
1239    37 http://www.sequenceontology.org/
1240    38 http://uberon.github.io/
1241    39 http://www.wormbase.org/
1242    40 http://wwf.panda.org/about_our_earth/ecoregions/ecoregion_list/
1243

1244    Table 2: Mean annual temperature and precipitation associated with miniature and non-miniature

1245    phenotypes in fishes

1246

| Variable | Type | Mean | p value | df | t Statistic | t Critical |
|---|---|---|---|---|---|---|
| Annual Mean Temperature (°C) | Miniature | 24.8°C | 0.002 | 227 | 3.128 | 1.970 |
| | Non-miniature | 22.6°C | | | | |
| Annual Mean Precipitation (mm) | Miniature | $6.9 \times 10^7$ mm | 0.008 | 227 | 2.668 | 1.970 |
| | Non-miniature | $1.8 \times 10^7$ mm | | | | |

1247

1248

1249    Table 3: Some of the URIs used to describe amphibian breeding and development in TraitBank

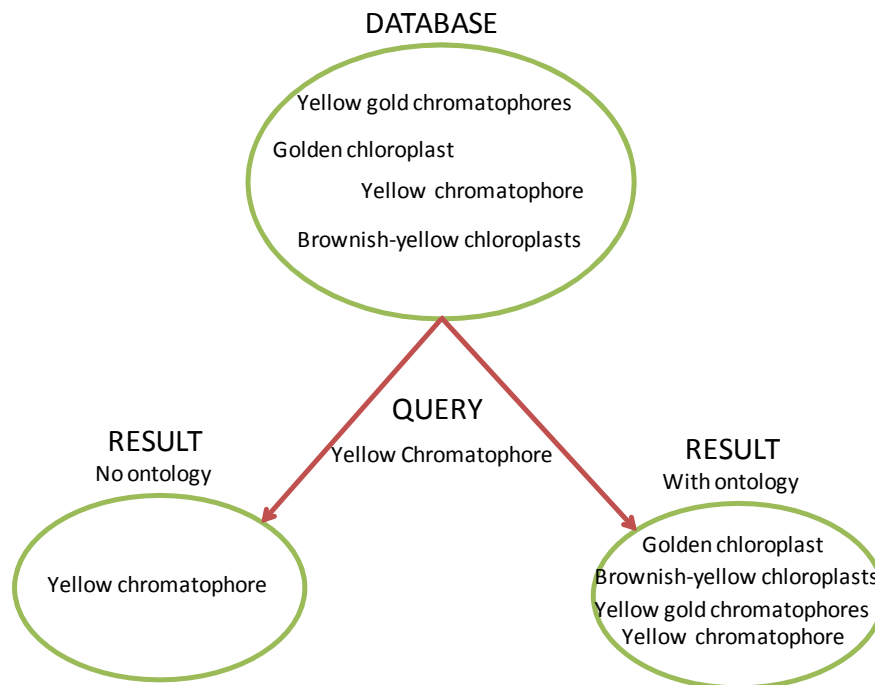| Term | URI |
| --- | --- |
| breeding habitat | http://eol.org/schema/terms/BreedingHabitat |
| development mode | http://eol.org/schema/terms/DevelopmentalMode |
| terrestrial habitat | http://purl.obolibrary.org/obo/ENVO_00002009 |
| intermittent pond | http://purl.obolibrary.org/obo/ENVO_00000504 |
| permanent pond | http://eol.org/schema/terms/permanentFreshwater |
| freshwater stream | http://eol.org/schema/terms/freshwaterStream |
| direct development | http://eol.org/schema/terms/directDeveloper |
| larval development | http://eol.org/schema/terms/larvalDevelopment |
| paedomorphic | http://purl.obolibrary.org/obo/HOM_0000029 |

1250

1251    Table 4: Breeding habitat and developmental mode for 282 species of amphibians

1252

|  | Larval | Direct | Paedomorphic | df | test statistic | $X^2_{0.95}$ |
|---|---|---|---|---|---|---|
| Freshwater Stream | 30 | 0 | 1 | 6 | 278 | 1.635 |
| Intermittent Pond | 28 | 0 | 0 |  |  |  |
| Permanent Pond | 59 | 2 | 3 |  |  |  |
| Terrestrial | 2 | 166 | 0 |  |  |  |

1253

1254     Figure 1: Ontology and the Heterogeneity Challenge



DATABASE

Yellow gold chromatophores

Golden chloroplast

Yellow chromatophore

Brownish-yellow chloroplasts

QUERY
Yellow Chromatophore

RESULT
No ontology

Yellow chromatophore

RESULT
With ontology

Golden chloroplast
Brownish-yellow chloroplasts
Yellow gold chromatophores
Yellow chromatophore

1255
1256     This diagram demonstrates how ontologies can solve the challenge of heterogeneous
1257     terminology. In this example, the database contains four different natural language descriptions
1258     about dinoflagellate chloroplasts harvested from text. A user needs to query the database for
1259     instances of dinoflagellates with yellow chromatophores. Without an ontology to provide the
1260     query engine information about synonomy ("chromatophore" = "chloroplast") and term
1261     relationships ("brownish-yellow", "golden", and "yellow gold" are subtypes of "yellow"), a
1262     query for "yellow chromatophore" will only yield one of the four results the user needs and
1263     would find using an ontology. Without an ontology to link closely related concepts with a
1264     common parent, and reconcile heterogeneous terms, a user would have to perform many more
1265     queries to get a desired result, which may not be tractable in a large dataset.
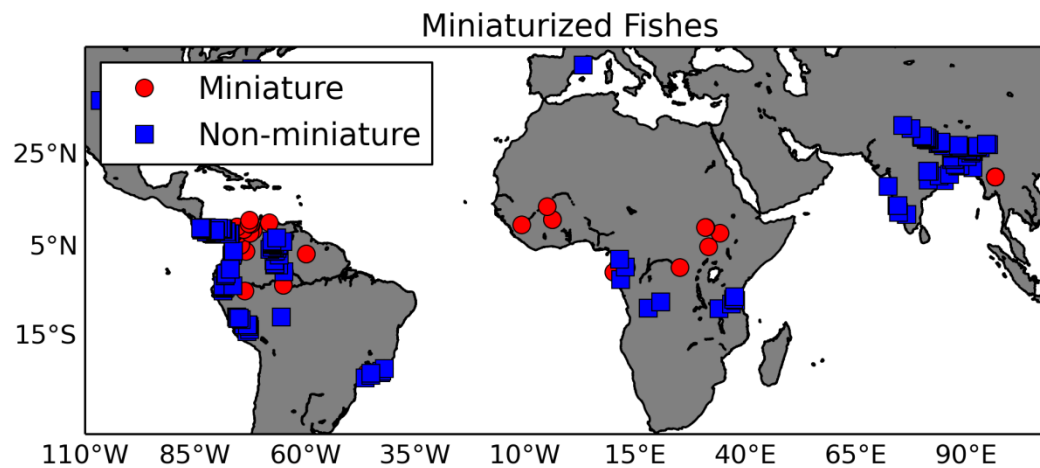
1266     Figure 2: Manual Workflow Conceptual Diagram

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280



1281    This diagram shows the manual workflow to link phenotype and environment data sets using

1282    current tools and services.

1283

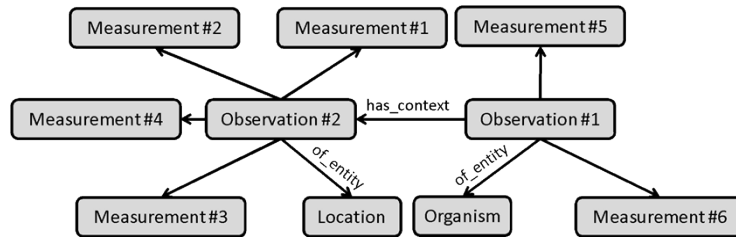1284     Figure 3: Map of Miniaturized Fishes and Their Non-Miniaturized Sister Taxa



1285

1286     This map shows locations of fish species exhibiting the miniaturized phenotype (red circles) and

1287     their non-miniature sister taxa (blue squares). The georeferenced occurrence data were gathered
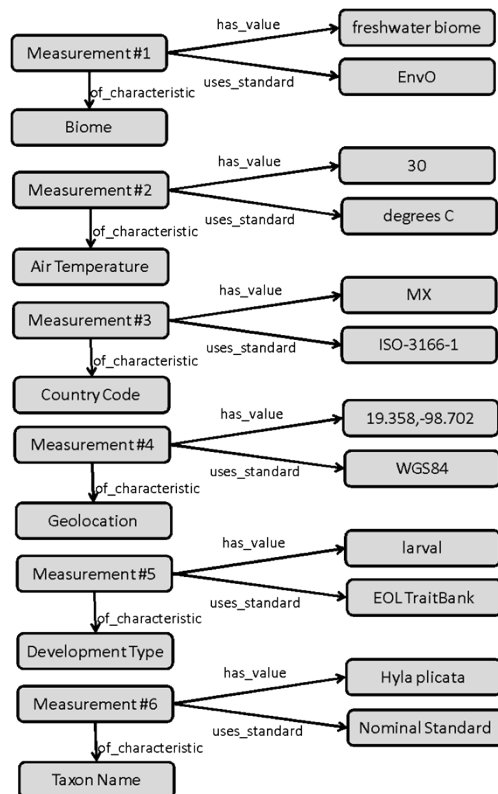
1288     from GBIF.

1289

1290    Figure 4: Using the OBOE Ontology to Link Phenotype and Environment

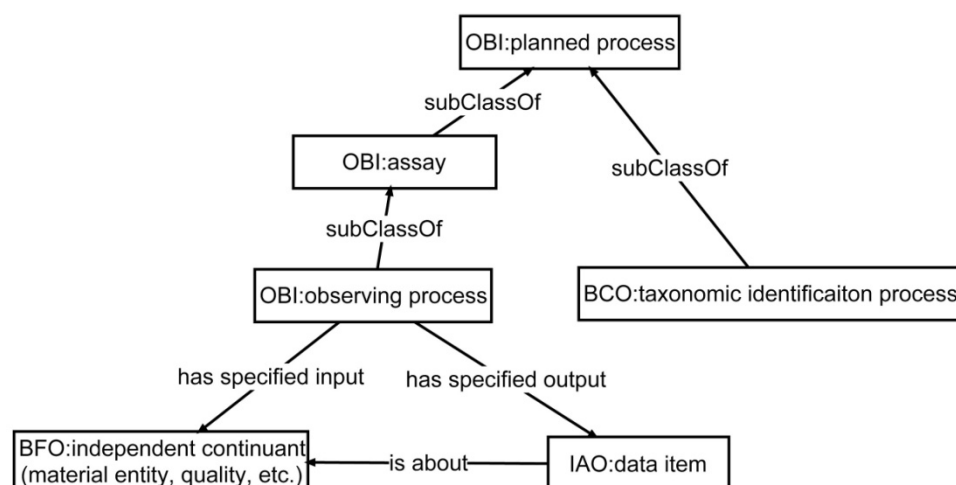1291    A



1292

1293    B



1294

1295    This demonstrates linking phenotype and environment using instances of the OBOE classes

1296    Entity, Observation, and Measurement. A) Links between Entity, Observation, and Measurement

1297    OBOE classes. B) Example measurements of phenotypes and environments using instances of

1298    the OBOE classes. Numbered measurement instances are consistent across A and B. This

1299    representation is simplified with regards to the taxonomic entities in play (Baskauf and Webb

1300    unpublished data).
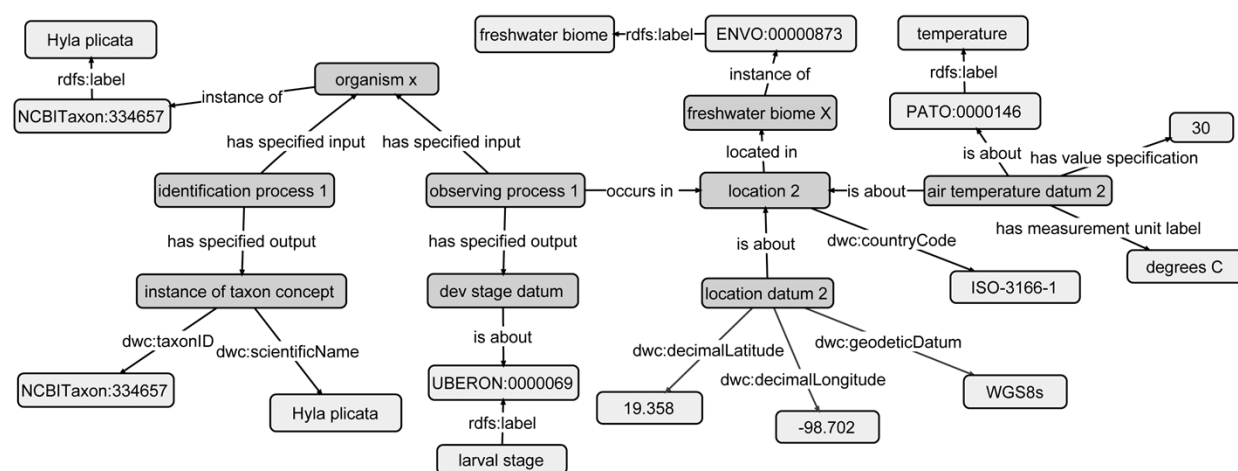
1301    Figure 5: Using BCO to Link Phenotype and Environment

1302    A



1303

1304    B



1305

1306    This demonstrates linking phenotype and environment using classes and relations from the

1307    Biological Collections Ontology (BCO). A) A simple version of the classes and relations used to

1308    describe observations in the BCO, with classes imported from OBI (Ontology for Biomedical

1309    Investigations), IAO (Information Artifact Ontology), and BFO (Basic Formal Ontology). B)

1310    Links among organism, phenotype, and environment, using the BCO model, using the same data

1311    as in Fig. 3. Light grey boxes represent either literal values (e.g., Hyla plicata), or instances of

1312    classes from external ontologies (ENVO – Environment Ontology, UBERON – Uber Anatomy

1313    Ontology, PATO – Phenotye Quality Ontology). Properties with a dwc prefix are imported

1314    directly from Darwin Core.