# Phenoscape: Semantic analysis of organismal traits and genes yields insights in evolutionary biology

Paula M. Mabee[1], Wasila M. Dahdul[1], James P. Balhoff [2], Hilmar Lapp [3], Prashanti Manda[4], Josef Uyeda[5], Todd Vision[6], Monte Westerfield[7]

[1]Department of Biology, University of South Dakota, Vermillion, South Dakota, USA
[2]Renaissance Computing Institute, University of North Carolina, Chapel Hill, North Carolina, USA
[3]Center for Genomic and Computational Biology, Duke University, Durham, North Carolina, USA
[4]Department of Computer Science, University of North Carolina at Greensboro, North Carolina, USA
[5]Department of Biological Sciences, Virginia Tech, Blacksburg, Virginia, USA
[6]Department of Biology, University of North Carolina at Chapel Hill, North Carolina, USA
[7]Institute of Neuroscience, University of Oregon, Eugene, Oregon, USA

Corresponding Author:
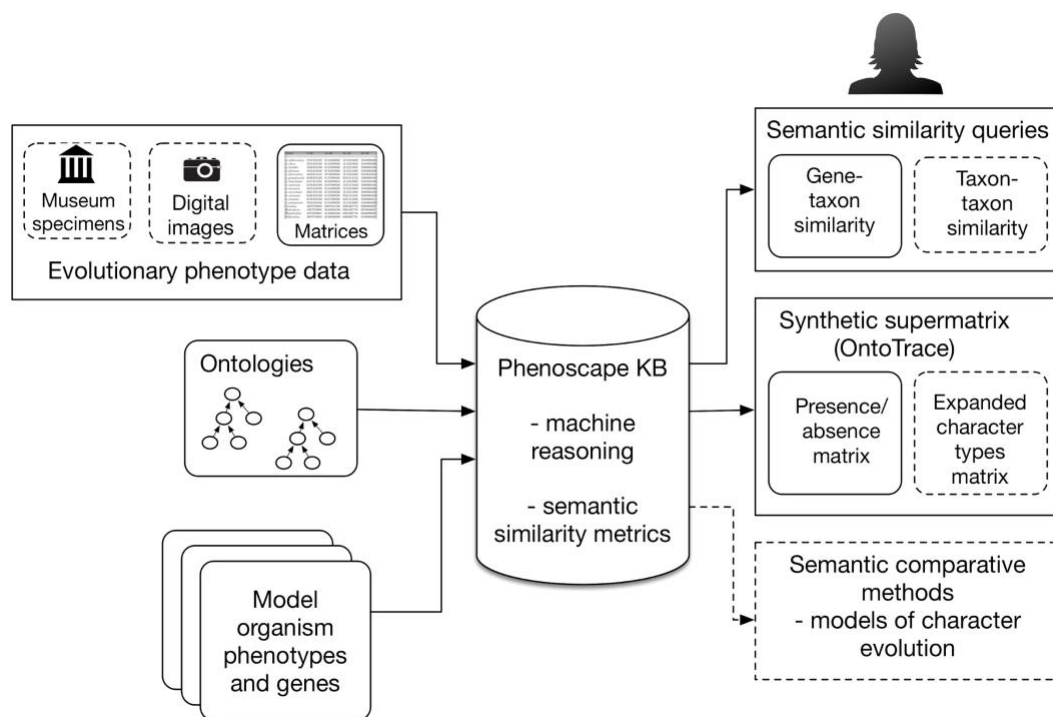Paula M. Mabee[1]
Email address: paula.mabee@usd.edu

# ABSTRACT

The study of how the observable features of organisms, i.e., their phenotypes, result from the complex interplay between genetics, development, and the environment, is central to much research in biology. The varied language used in the description of phenotypes, however, impedes the large scale and interdisciplinary analysis of phenotypes by computational methods. The Phenoscape project (www.phenoscape.org) has developed semantic annotation tools and a gene–phenotype knowledgebase, the Phenoscape KB, that uses machine reasoning to connect evolutionary phenotypes from the comparative literature to mutant phenotypes from model organisms. The semantically annotated data enables the linking of novel species phenotypes with candidate genes that may underlie them. Semantic annotation of evolutionary phenotypes further enables previously difficult or novel analyses of comparative anatomy and evolution. These include generating large, synthetic character matrices of presence/absence phenotypes based on inference, and searching for taxa and genes with similar variation profiles using semantic similarity. Phenoscape is further extending these tools to enable users to automatically generate synthetic supermatrices for diverse character types, and use the domain knowledge encoded in ontologies for evolutionary trait analysis. Curating the annotated phenotypes necessary for this research requires significant human curator effort, although semi-automated natural language processing tools promise to expedite the curation of free text. As semantic tools and methods are developed for the biodiversity sciences, new insights from the increasingly connected stores of interoperable phenotypic and genetic data are anticipated.

**INTRODUCTION**

21

22        There are over 20 million extant species on the planet, most of which can be described in

23   relation to their unique and widely diverse phenotypes. Comparisons across species phenotypes,

24   however, cannot yet readily be made using computer-assisted methods. This is because the rich

25   legacy of comparative morphology has not yet been semantically enabled—that is, the corpus is

26   in a free-text format that renders computation nearly impossible. This situation began to change

27   almost two decades ago when model organism geneticists began representing the phenotypic

28   changes resulting from experimental gene manipulations, with terms from anatomy or phenotype

29   ontologies that they developed for each model organism (e.g., Sprague et al. 2001).  More

30   recently, the opportunity to enable interoperability from the phenotypes of biodiverse species to

31   candidate genes from model species (Mabee et al. 2007a, 2007b) motivated the Phenoscape team

32   to develop one of the first multispecies anatomy ontologies, the Teleost Anatomy Ontology

33   (Dahdul et al. 2010b), based initially on the Zebrafish Anatomy Ontology (Ruzicka et al. 2015).

34   Developing ontologies appropriate for biodiversity, including taxonomy ontologies (Midford et

35   al. 2013) and scaling them up first to the level of teleost fishes (Dahdul et al. 2010), then to the

36   level of vertebrates (Dahdul et al. 2012) and then to the level of metazoans (Mungall et al. 2012;

37   Haendel et al. 2014), further enabled the automation of phenotypic comparisons across vertebrate

38   species and discovery of candidate genes underlying evolutionarily novel phenotypes by the

39   team (Edmunds et al. 2016).  Over the past ten years a broad community of scientists invested in

40   the development of shared community ontologies (e.g., Gkoutos et al. 2005; Haendel et al. 2008,

41   2014; Dahdul et al. 2014), annotation tools (Balhoff et al. 2010, 2014a; Yoder et al. 2010; Cui et

42   al. 2016; The Gene Ontology Consortium 2017) and formats (Dahdul et al. 2010a; Vos et al.

43   2012) for phenotype annotation across biodiverse species (Dahdul et al. 2010a). These resources

2

44    have made computational analyses possible and they have been leveraged to build a wealth of

45    innovative applications (e.g., Deans et al. 2012; Mullins et al. 2012; Balhoff et al. 2013;

46    Dececchi et al. 2015; Manda et al. 2015; Druzinsky et al. 2016; Jackson et al. 2018) across a

47    variety of biodiversity-based research. The Phenoscape Knowledgebase (KB) (Figure 1)

48    demonstrates these connections by integrating gene phenotype annotations from model organism

49    databases with phenotype annotations from the biodiversity literature (Table 1). Compelling

50    demonstrations of the utility of semantics for biodiversity studies are important because of the

51    large and expensive investments in infrastructure and tool development required to curate the

52    legacy literature and move the publication of phenotypic data into a natively semantic form.



53
54    **Figure 1.** Flow chart of currently existing data sources and tools (solid borders and lines) in the

55    Phenoscape KB, and data and tools not yet integrated or developed (dotted borders and lines) but

56    relevant to users in biodiversity research.

57

3

58    **Table 1:** Data for evolutionary and model organism phenotypes in the Phenoscape KB. (Data as

59    of 2018-05-11)

60    **Evolutionary Phenotypes**

| | |
|---|---|
| Annotated anatomical character states | 22,321 |
| Total number of annotated taxa (extant and fossil vertebrates) | 5,310 |
| Total number of taxon phenotypes | 540,163 |
| Terminal taxa (species) with at least one phenotype | 4,260 |
| Non-terminal taxa with at least one phenotype | 1,050 |
| Evolutionary phenotype profiles | 682 |

61

62    **Model Organism Phenotypes**

| | Zebrafish | Mouse | Xenopus | Human |
|---|---|---|---|---|
| Genes with at least one phenotype | 5,883 | 7,758 | 12 | 3,717 |
| Phenotype annotations | 90,132 | 171,876 | 236 | 123,956 |
| Genes with any expression data | 12,509 | 10,599 | 15,062 | 0 |
| Gene expression annotations | 179,232 | 800,824 | 454,337 | 0 |

63      To date, only a small proportion of the biodiversity literature has been annotated

64   semantically, and no publisher, to our knowledge, tags phenotypes with ontological terms that

65   would support interoperability. The comparative study of organismal phenotypes, however,

66   motivates research across diverse fields of biology, including evolution, paleontology,

67   developmental biology, agriculture, and the veterinary and health sciences (Deans et al. 2015).

68   The efficiency and potential of fundamental discoveries in the biodiversity arena would be

69   dramatically expanded by the increased use of semantics. Further, few species, i.e., only model

70   organisms, have curated phenotypic data that is linked to genetic and genomic data. The growth

71   in sequencing technology, however, is changing this dynamic, resulting in the rapid expansion of

72   genomic data for non-model species (e.g., Russell et al. 2017 and Chapter 10). However, without

73   corresponding phenomic databases, the challenge of relating the growing volume of genetic

74   knowledge in model and emerging model organisms to the diversity of phenotypes in nature

75   cannot be met.  In this chapter, through the description of driving research questions and by

76   examples of the use of semantically annotated data in the Phenoscape KB, we provide a glimpse

77   of the promise that semantic analysis tools hold in comparing phenotypes across species and

78   globally associating genetic to phenotypic data.

79   **1. Relating biodiverse phenotypes to candidate genes**

80      Identifying the genetic and developmental changes that brought forth the incredible

81   phenotypic diversification of life is a recalcitrant problem, but one where a basic semantic

82   approach has shown promise and where more sophisticated approaches using semantic similarity

83   may yet be even more valuable. Semantic similarity enables comparison and analysis of semantic

84   annotations between entities (genes, taxa) using ontologies and computational reasoners to

85   compute scores that reflect the level of similarity (e.g., Washington et al. 2009; Manda et al.

86  2015; see examples in Chapter 10). The Phenoscape team showed that ontology-driven

87  information systems can generate thousands of testable hypotheses relating unique morphologies

88  from non-model biodiverse species to candidate genes (Mabee et al. 2012).  One of these, for

89  example, connected the unique loss of a tongue ('basihyal element') in catfishes (Siluriformes)

90  with several candidate genes from the zebrafish data.  Edmunds et al. (2016) experimentally

91  tested the candidates by examining their endogenous expression patterns in the channel catfish,

92  *Ictalurus punctatus*, and found results consistent with the in silico hypothesis that the tongue

93  evolved through disruption in developmental pathways at, or upstream of, *brpf1*.

94         The Phenoscape team recently extended this approach (Manda et al. 2015) by using

95  semantic similarity to find matches between the full set of phenotypes described for a gene and

96  the unique set of phenotypes that characterizes a clade of species, i.e., an 'evolutionary

97  phenotype profile'. The effects from a gene knockdown range from several to hundreds of

98  phenotypes, and the goal is to compare these in their entirety to the calculated set of phenotypes

99  that are variable among the immediate descendants of a particular taxon. Using semantic

100  similarity, the Phenoscape KB performs fuzzy matching between suites of phenotypes, and

101  displays the taxonomic groups that vary in phenotypes that match most closely to the gene

102  profile that results when the action of that gene is disrupted (e.g., knocked down). The user

103  interface provides the statistical support for each match and allows the supporting evidence to be

104  examined. There are some important caveats that must be considered when interpreting the

105  results, such as the potential for some matches to result from differences in annotation coverage

106  between genetic and evolutionary studies in the KB. Potentially spurious matches in that

107  category are flagged by the KB. The KB also provides an interface for the reverse query: what

108  genes have phenotypes that match most closely to the set of evolutionary phenotypes in a

6

109　　particular taxon under consideration? That is, a biologist who is curious about the genetic basis

110　　of taxonomic diversity might want to find genes that have phenotypes that resemble the

111　　phenotypic variation exhibited by a particular taxon.


112　　**2. Future applications of semantic similarity to phenotypes of biodiverse taxa**

113　　　　　Questions of whether a particular combination of phenotypes in a taxon is unique, or

114　　what it might be similar to, are the types of broad questions that may be addressed in applying

115　　semantic similarity-based data mining to phenotypes across diverse taxa. Semantic similarity

116　　would retrieve taxa with similar phenotypic profiles; such similarity may have arisen because of

117　　common ancestry or independent origin (a 'homoplasy finder'). As described by Braun et al.

118　　(Chapter 10), predictive phenomics can, for example, be used to target desired phenotypes in

119　　species of interest - and together with recent gene editing capabilities, functional genomic

120　　analysis can be newly brought to bear on biodiverse species. The Phenoscape KB currently

121　　enables users to view taxa with variation similar to the phenotypic profile of a gene (and *vice*

122　　*versa*). In the future, they will also be able to query one custom set of phenotypes against another

123　　or a taxonomically selected subset, and obtain a ranked list of taxa with similar phenotypes. For

124　　example, miniature fishes in the genus *Paedocypris*, like many fishes that are evolutionarily

125　　reduced to an extremely small body size, exhibit the absence of bones including the interhyal,

126　　vomer, parietal, posttemporal, and supraneurals (Britz and Conway 2009). Are there other taxa

127　　that lack a highly similar set of bones? Enabling a comparison of these phenotypes across

128　　diverse taxa would allow a user to query for such matches; in this case, matches would include

129　　the ricefishes in the family Adrianichthyidae (Wiley and Johnson 2010), which similarly lack the

130　　interhyal, vomer, and supraneurals, and other bones such as the supracleithrum. Further,

131　　adrianichthyids may lack or possess extremely small or absent parietal bones and have

7

132   structurally simple posttemporal bones, which biologists may recognize as reductive phenotypes

133   on a continuum close to 'absent'.  Methods that incorporate a framework of probabilistic

134   reasoning for phenotype relatedness (e.g., Bauer et al. 2012) have the potential to improve

135   precision of ontology-based queries.


136   **3.   Relating biodiverse phenotypes across studies: presence/absence**

137          Addressing many of the questions in the biodiversity sciences involve knowing how a

138   specific trait or set of traits has evolved across a group of species.  Although the published

139   literature is replete with research relating species and traits, and a few repositories hold

140   phylogenetic trees, some of which are computed products from trait data, neither the traits nor

141   the trees can be easily synthesized across studies. The OntoTrace tool was developed by the

142   Phenoscape team (Balhoff et al. 2014b; Dececchi et al. 2015) to enable users to automatically

143   pull together, from phenotype annotations made to published character matrices and

144   monographic texts (Dececchi et al. 2015, 2016), a set of presence/absence data for specific traits

145   for a set of taxa.  For example, querying the Phenoscape KB for a supermatrix of traits of fins,

146   limbs, girdles and their parts in sarcopterygian vertebrates (lobe-finned fishes and tetrapods),

147   Dececchi et al. (2015) retrieved data for 1,052 taxa from 55 studies. The data, 1,759 synthetic

148   presence/absence characters, were derived from 2,588 text-based character states (1,195

149   characters). The resultant character by taxon matrix was termed a 'synthetic morphological

150   supermatrix'. Because of the ontological annotations, not only could these phenotypic data be

151   automatically aggregated from multiple studies into a supermatrix, but the asserted data could be

152   extended through inference to traits that were implied by, but not directly asserted in the original

153   publications. For example, if an author observed a curved pectoral fin ray in a species, the

154   machine would infer, based on the knowledge of anatomy encoded in the requisite ontology

8

155    (Uberon in this case), that a pectoral fin is present in that species (see Dececchi et al. 2015 and

156    Jackson et al. 2018 for further examples). In this manner, the missing data in the variable

157    character subset of the matrix (the subset containing only characters that include both present and

158    absent states) was reduced from 98.5% to 78.2%. Further, 76% of the variable characters were

159    made variable through the addition of inferred states. The authors pointed out that character

160    conflicts and provenance reports from OntoTrace would support researchers review of large

161    aggregated data sets and they showed how such machine reasoning enables quantification and

162    new visualizations of the data, allowing the identification of undersampled character space.


163        **4.  Relating biodiverse phenotypes to phylogenetic trees**

164        Using ontologies and machine reasoning to automatically generate large, synthetic

165    character matrices of presence/absence phenotypes (as per above) set the stage for the research of

166    Jackson et al. (2018), who took this a step further.  They developed a bioinformatic pipeline to

167    propagate data that was asserted to higher-level taxonomic nodes, to descendant species that

168    were missing data.  Similar to Dececchi et al. (2015), they showed that such logic inference

169    significantly extended the asserted data (missing data were reduced from 98.0% to 85.9%), but

170    additionally they showed the value of taxonomic data propagation, which extended the data

171    further, reducing missing data to 34.8% (Jackson et al. 2018). Using the resultant matrix along

172    with a synthetic phylogeny from the Open Tree of Life (Hinchliff et al. 2015), they mapped the

173    full trait data set for 12,582 species to the tree and addressed the question of how often paired

174    fins were lost in teleost fishes and whether they were ever regained (Jackson et al. 2018).

175    Looking ahead, if all published traits and trees were made computable using these methods, any

176    user could automatically generate a matrix for a specified set of traits and map it on various

177    synthetic tree topologies, which in turn would allow addressing a host of questions regarding the

9

178     pattern and tempo of phenotypic evolution and associations with genomic and environmental

179     (Thessen et al. 2015) variables.


180     **5.  Relating biodiverse phenotypes across studies: future work**

181          As described above, OntoTrace generates synthetic morphological supermatrices for

182     presence/absence characters only (Dececchi et al. 2015).  Expanding this functionality to

183     automatically synthesize characters of other qualities, such as shape, size, structure, and color, is

184     a current challenge that the Phenoscape team is addressing. For example, whereas characters in a

185     presence/absence matrix are by definition limited to two states per character, the number of

186     possible states for characters in other categories is *a priori* unconstrained. Thus, automatically

187     synthesizing characters that, for example, describe 'basihyal bone, shape', can result in a large

188     number of states per character because every originally published state that semantically is some

189     type of 'basihyal bone shape' would have to be appended as a new state to the synthesized

190     character.  In the case of this example, there may be seven distinct shape terms used in its

191     annotation (Box 1). The ontological relationships indicate that subsets of these states are more

192     similar to each other than others. By adapting current semantic similarity metrics for the purpose

193     of character and character state aggregation, and in effect, homology assignment, these distinct

194     shape descriptors can be consolidated into new, synthetic states (see matrix in Box 1).

195

196 ⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯

197 **Box 1. Assembling a synthetic character and its states for 'basihyal bone, shape'.**

198 *Step 1: Assemble list of 'shape' (PATO:0000052) quality terms for all characters and states from*

199 *multiple publications that include the entity 'basihyal bone' (UBERON:0011618):*

200 'spiny' (PATO:0001365)

201 'folded' (PATO:0001910)

202 'upturned' (PATO:0002031)

203 'blade-like' (PATO:0002235)

204 'pointed' (PATO:0002258)

205 'curved ventral' (PATO:0001469)

206 'tapered' (PATO:0001500)

207

208 *Step 2: Apply semantic similarity to above list of PATO terms for basihyal bone.  Because of*

209 *higher similarity among terms, three states (0, 1, 2) are generated from the seven phenotypes:*

210 **Character 1: Basihyal bone: shape**

211 Synthetic State 0: 'sharp' (PATO:0001419) (includes 'blade-like', 'pointed', 'tapered')

212 Synthetic State 1: 'curved' (PATO:0000406) (includes 'upturned', 'curved ventral')

213 Synthetic State 2: 'surface feature shape' (PATO:0001925) (includes 'spiny', 'folded')

214 ⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯

215

216          The Phenoscape team is now developing semantic similarity-based methods to cluster

217 phenotypes across different character categories into characters and states, thus automating

218 matrix construction, and enabling users to optimize the matrix for a variety of metrics. This

219   would allow a user to constrain the number of characters in a synthesized matrix by excluding

220   those with low information content (e.g., those for high level terms from the anatomy ontology

221   such as 'fin'' vs. 'pectoral fin'). Thus, employing semantic reasoning in matrix construction will

222   allow a user to balance the properties of a synthetic matrix between, on the one hand, containing

223   highly specific characters (and thus increased missing data), and on the other, including lower

224   specificity characters (and thus decreasing missing data).

225        In addition to semantic tools for supermatrix construction, the Phenoscape team is

226   developing enhanced semantics for addressing questions of trait evolution. Unlike the current

227   tools available for analyzing molecular data, where each nucleotide site can be treated as

228   independent of each other, evolutionary models for large morphological character matrices face

229   significant challenges overcoming the strong conditional dependencies and correlations among

230   morphological traits. Most existing methods ignore such dependencies and morphological

231   characters are treated as independent. By leveraging domain knowledge relevant to assessing

232   correlations of the traits underlying the characters, Phenoscape is developing tools that enable

233   users to incorporate evidence of the relatedness of traits in a morphological matrix and into

234   models of character evolution. These include measures of trait independence based on

235   ontological relationships, distance (semantic similarity) of traits in the knowledge graph, and

236   measures of genetic overlap (as derived from gene-phenotype annotations from model organism

237   databases). Such dependencies can be directly built into the macroevolutionary model, or can be

238   used to inform prior probabilities in Bayesian analyses when grouping traits into modules with

239   shared evolutionary parameters or dynamics.

240        One of the challenges in conducting semantic similarity comparisons is the computational

241   overhead of comparing EQ phenotypes over a large ontology space. Improvements in scalability

242   of semantic similarity methods would enable fast, on-the-fly semantic similarity searches.

243   Successfully applying these methods also currently depends on accurate *a posteriori* annotation

244   of characters to capture the original author's intent.  With only the published description and

245   perhaps images to rely on, curators are unable to consistently apply standardized terms, a factor

246   leading to lower consistency (Cui et al. 2015).  For example, in a comparison of curator vs.

247   machine generated annotations (Cui et al. 2015; in prep), three curators described the increased

248   distance between the contralateral pelvic fins with three different quality descriptors: 'far from',

249   'separated from', and 'set apart from'.  As methods and software tools develop, such that original

250   authors are better empowered to apply the semantics themselves, the accuracy of character

251   annotation, and thus, consolidation will increase. In the above example, the author would

252   presumably be able to choose, based on the different definitions of the ontology terms, which

253   term is most applicable to the phenotype observed.


254   **6.  Future challenges**

255       A long-standing question, and one also being currently tackled by the Phenoscape team,

256   is how the relationship of historical homology, i.e., similarity due to common ancestry, can most

257   effectively be used in data retrieval.  Recent work by Manda et al. (2016b), examined how

258   semantic similarity is affected when external homology knowledge is included in an ontology.

259   They measured phenotypic similarity between orthologous and non-orthologous gene pairs

260   between humans and either mouse or zebrafish, and they compared the effect of including real

261   vs. faux homology axioms. Semantic similarity was preferentially increased for orthologs when

262   using real homology axioms, though only across the more divergent of the two species (human to

263   zebrash, not human to mouse) (Manda et al., 2016).  Overall, the effect of including homology

264   axioms on cross-species semantic similarity was modest, though the authors suggested that the

265   effect might be greater for more distant species comparisons. Current efforts include editing and

266   clarifying the homology relationships in the Uberon ontology and investigating how reasoning

267   on different models of homology affects information retrieval in the KB.

268         Another challenge for the broader application of semantics to biodiversity data is the

269   significant, largely manual, effort necessary to annotate phenotypes from the published literature

270   (Dahdul et al. 2015). Natural language processing tools are needed going forward to auto-

271   annotate the legacy literature (Arighi et al. 2013; Cui et al. 2015; in prep). Further, in the future

272   semantic phenotype data may increasingly come directly from publications, as semi-automated

273   methods for marking up manuscripts at the time of publication become more accurate, mature,

274   and thus prevalent. Evaluating, and hence continuously improving the accuracy of machine

275   generated annotations depends on expert-curated "gold standard" data sets. To this end,

276   Phenoscape has developed the first gold standard dataset for biodiversity phenotypes (in prep).

277   Efforts to use ontologies in the process of new species descriptions are underway (Deans et al.

278   2012; Balhoff et al. 2013), and will contribute to achieving a vision of widely available  linked

279   species phenotype data.

280         As high-throughput phenotyping, typically involving image data collection, becomes

281   more scalable, the application of semantic metadata would enable automated connections to the

282   tools and computable datasets described herein. These digitization efforts can be new sources of

283   phenotype information (Figure 1). Although broad domains of biology can be served if semantics

284   are placed on digitized images and specimens, so far only a few projects are using semantics to

285   label digitized specimens and their parts, despite promising prototypes (Maglia et al. 2007;

286   Rámirez et al. 2007).  If anatomical parts were tagged with ontology terms, then queries on basic

287   trait distributions could be enabled (e.g., presence of pectoral fins in taxa a, b, c...). Although

288    having a reduced information content compared to full Entity-Quality expressions, entity-only

289    annotations have been shown to be informative for semantic similarity (Manda et al. 2016a).

290    Thus, new sources of phenotypic data, such as those for specimens of extinct and extant taxa

291    associated with institutional collections, can easily be made interoperable through shared

292    semantics (Figure 1).

293    **CONCLUSIONS**

294    Over the past 10 years the development of shared cross-species community ontology resources

295    such as Uberon and PATO has enabled interoperability of phenotype and genotype data.  This in

296    turn enables a wealth of potential applications and discoveries from semantic analysis of

297    biodiverse taxa.  Scientific attention continues to move toward gaining a deeper fundamental

298    understanding of the developmental and evolutionary relationship between genotype and

299    phenotype.  The profound scale and scope of this problem will not only require interoperable big

300    data, both genomic and phenomic, from a biodiverse set of taxa, but also new ways of using

301    machines to enable this understanding.  The applications of semantic analysis described herein

302    only scratch the surface of what is possible. As scientific publication moves to incorporate

303    semantic markup of phenotype data, and semi-automated tools are improved to annotate the

304    phenotype legacy literature, knowledge of the rich phenotypic palette of life on our planet can be

305    exposed to machine computation with great advantage to fundamental discovery across the life

306    sciences.

316 **REFERENCES**

317 Arighi C.N., Carterette B., Cohen K.B., Krallinger M., Wilbur W.J., Fey P., Dodson R., Cooper

318     L., Van Slyke C.E., Dahdul W., Mabee P., Li D., Harris B., Gillespie M., Jimenez S.,

319     Roberts P., Matthews L., Becker K., Drabkin H., Bello S., Licata L., Chatr-aryamontri A.,

320     Schaeffer M.L., Park J., Haendel M., Van Auken K., Li Y., Chan J., Muller H.-M., Cui H.,

321     Balhoff J.P., Chi-Yang Wu J., Lu Z., Wei C.-H., Tudor C.O., Raja K., Subramani S.,

322     Natarajan J., Cejuela J.M., Dubey P., Wu C. 2013. An overview of the BioCreative 2012

323     Workshop Track III: interactive text mining task. Database. 2013:bas056.

324 Balhoff J.P., Dahdul W.M., Dececchi T.A., Lapp H., Mabee P.M., Vision T.J. 2014a. Annotation

325     of phenotypic diversity: decoupling data curation and ontology curation using Phenex. J.

326     Biomed. Semantics. 5:45.

327 Balhoff J.P., Dahdul W.M., Kothari C.R., Lapp H., Lundberg J.G., Mabee P., Midford P.E.,

328     Westerfield M., Vision T.J. 2010. Phenex: ontological annotation of phenotypic diversity.

329     PLoS One. 5(5):e10500.

330 Balhoff J.P., Dececchi T.A., Mabee P.M., Lapp H. 2014b. Presence-absence reasoning for

331    evolutionary phenotypes. Peer-reviewed conference paper, Bio-ontologies SIG at ISMB

332    2014.

333    Balhoff J.P., Mikó I., Yoder M.J., Mullins P.L., Deans A.R. 2013. A semantic model for species

334    description applied to the ensign wasps (Hymenoptera: Evaniidae) of New Caledonia. Syst.

335    Biol. 62:639–659.

336    Bauer S., Köhler S., Schulz M.H., Robinson P.N. 2012. Bayesian ontology querying for accurate

337    and noise-tolerant semantic searches. Bioinformatics. 28:2502–2508.

338    Britz R., Conway K.W. 2009. Osteology of Paedocypris, a miniature and highly developmentally

339    truncated fish (Teleostei: Ostariophysi: Cyprinidae). J. Morphol. 270:389–412.

340    Cui H., Dahdul W., Dececchi T.A., Ibrahim N., Mabee P., Balhoff J.P., Gopalakrishnan H. 2015.

341    CharaParser+EQ: Performance Evaluation without Gold Standard. Proceedings of the 78th

342    Annual Meeting of the Association for Information Science and Technology Annual

343    (ASIS&T). Vol. 51.

344    Cui H., Xu D., Chong S.S., Ramirez M., Rodenhausen T., Macklin J.A., Ludäscher B., Morris

345    R.A., Soto E.M., Koch N.M. 2016. Introducing Explorer of Taxon Concepts with a case

346    study on spider measurement matrix building. BMC Bioinformatics. 17:471.

347    Dahdul W., Dececchi T.A., Ibrahim N., Lapp H., Mabee P. 2015. Moving the mountain: analysis

348    of the effort required to transform comparative anatomy into computable anatomy. Database

349    . 2015:bav040.

350    Dahdul W.M., Balhoff J.P., Blackburn D.C., Diehl A.D., Haendel M.A., Hall B.K., Lapp H.,

351    Lundberg J.G., Mungall C.J., Ringwald M., Segerdell E., Van Slyke C.E., Vickaryous M.K.,

352    Westerfield M., Mabee P.M. 2012. A unified anatomy ontology of the vertebrate skeletal

353        system. PLoS One. 7(12):e51070.

354    Dahdul W.M., Balhoff J.P., Engeman J., Grande T., Hilton E.J., Kothari C., Lapp H., Lundberg

355        J.G., Midford P.E., Vision T.J., Westerfield M., Mabee P.M. 2010a. Evolutionary

356        characters, phenotypes and ontologies: curating data from the systematic biology literature.

357        PLoS One. 5(5):e10708.

358    Dahdul W.M., Cui H., Mabee P.M., Mungall C.J., Osumi-Sutherland D., Walls R.L., Haendel

359        M.A. 2014. Nose to tail, roots to shoots: spatial descriptors for phenotypic diversity in the

360        Biological Spatial Ontology. J. Biomed. Semantics. 5:34.

361    Dahdul W.M., Lundberg J.G., Midford P.E., Balhoff J.P., Lapp H., Vision T.J., Haendel M.A.,

362        Westerfield M., Mabee P.M. 2010b. The Teleost Anatomy Ontology: anatomical

363        representation for the genomics age. Syst. Biol. 59:369–383.

364    Deans A.R., Lewis S.E., Huala E., Anzaldo S.S., Ashburner M., Balhoff J.P., Blackburn D.C.,

365        Blake J.A., Burleigh J.G., Chanet B., Cooper L.D., Courtot M., Csösz S., Cui H., Dahdul

366        W., Das S., Dececchi T.A., Dettai A., Diogo R., Druzinsky R.E., Dumontier M., Franz

367        N.M., Friedrich F., Gkoutos G.V., Haendel M., Harmon L.J., Hayamizu T.F., He Y., Hines

368        H.M., Ibrahim N., Jackson L.M., Jaiswal P., James-Zorn C., Köhler S., Lecointre G., Lapp

369        H., Lawrence C.J., Le Novère N., Lundberg J.G., Macklin J., Mast A.R., Midford P.E.,

370        Mikó I., Mungall C.J., Oellrich A., Osumi-Sutherland D., Parkinson H., Ramírez M.J.,

371        Richter S., Robinson P.N., Ruttenberg A., Schulz K.S., Segerdell E., Seltmann K.C.,

372        Sharkey M.J., Smith A.D., Smith B., Specht C.D., Squires R.B., Thacker R.W., Thessen A.,

373        Fernandez-Triana J., Vihinen M., Vize P.D., Vogt L., Wall C.E., Walls R.L., Westerfield M.,

374    Wharton R.A., Wirkner C.S., Woolley J.B., Yoder M.J., Zorn A.M., Mabee P. 2015. Finding

375    our way through phenotypes. PLoS Biol. 13(1):e1002033.

376    Deans A.R., Yoder M.J., Balhoff J.P. 2012. Time to change how we describe biodiversity.

377    Trends Ecol. Evol. 27:78–84.

378    Dececchi T.A., Balhoff J.P., Lapp H., Mabee P.M. 2015. Toward synthesizing our knowledge of

379    morphology: using ontologies and machine reasoning to extract presence/absence

380    evolutionary phenotypes across studies. Syst. Biol. 64:936–952.

381    Dececchi T.A., Mabee P.M., Blackburn D.C. 2016. Data sources for trait databases: comparing

382    the phenomic content of monographs and evolutionary matrices. PLoS One.

383    11(5):e0155680.

384    Druzinsky R.E., Balhoff J.P., Crompton A.W., Done J., German R.Z., Haendel M.A., Herrel A.,

385    Herring S.W., Lapp H., Mabee P.M., Muller H.-M., Mungall C.J., Sternberg P.W., Van

386    Auken K., Vinyard C.J., Williams S.H., Wall C.E. 2016. Muscle Logic: new knowledge

387    resource for anatomy enables comprehensive searches of the literature on the feeding

388    muscles of mammals. PLoS One. 11(2):e0149102.

389    Edmunds R.C., Su B., Balhoff J.P., Eames B.F., Dahdul W.M., Lapp H., Lundberg J.G., Vision

390    T.J., Dunham R.A., Mabee P.M., Others. 2016. Phenoscape: identifying candidate genes for

391    evolutionary phenotypes. Mol. Biol. Evol. 33:13–24.

392    Gkoutos G.V., Green E.C.J., Mallon A.-M., Hancock J.M., Davidson D. 2005. Using ontologies

393    to describe mouse phenotypes. Genome Biol. 6:R8.

394    Haendel M.A., Balhoff J.P., Bastian F.B., Blackburn D.C., Blake J.A., Bradford Y., Comte A.,

19

395    Dahdul W.M., Dececchi T.A., Druzinsky R.E., Hayamizu T.F., Ibrahim N., Lewis S.E.,

396    Mabee P.M., Niknejad A., Robinson-Rechavi M., Sereno P.C., Mungall C.J. 2014.

397    Unification of multi-species vertebrate anatomy ontologies for comparative biology in

398    Uberon. J. Biomed. Semantics. 5:21.

399    Haendel M.A., Neuhaus F., Osumi-Sutherland D., Mabee P.M., Mejino J.L.V., Mungall C.J.,

400    Smith B. 2008. CARO--the Common Anatomy Reference Ontology. In: Burger A,

401    Davidson D, Baldock R (eds) Anatomy Ontologies for Bioinformatics. Computational

402    Biology, vol 6. Springer, London.

403    Hinchliff C.E., Smith S.A., Allman J.F., Burleigh J.G., Chaudhary R., Coghill L.M., Crandall

404    K.A., Deng J., Drew B.T., Gazis R., Gude K., Hibbett D.S., Katz L.A., Laughinghouse H.D.

405    4th, McTavish E.J., Midford P.E., Owen C.L., Ree R.H., Rees J.A., Soltis D.E., Williams T.,

406    Cranston K.A. 2015. Synthesis of phylogeny and taxonomy into a comprehensive tree of

407    life. Proc. Natl. Acad. Sci. U. S. A.

408    Jackson L.M., Fernando P.C., Hanscom J.S., Balhoff J.P., Mabee P.M. 2018. Automated

409    integration of trees and traits: a case study using paired fin loss across teleost fishes. Syst.

410    Biol.:syx098.

411    Mabee P., Balhoff J.P., Dahdul W.M., Lapp H., Midford P.E., Vision T.J., Westerfield M. 2012.

412    500,000 fish phenotypes: The new informatics landscape for evolutionary and

413    developmental biology of the vertebrate skeleton. J. Appl. Ichthyol. 28:300–305.

414    Mabee P.M., Arratia G., Coburn M., Haendel M., Hilton E.J., Lundberg J.G., Mayden R.L., Rios

415    N., Westerfield M. 2007a. Connecting evolutionary morphology to genomics using

416     ontologies: a case study from Cypriniformes including zebrafish. J. Exp. Zool. B Mol. Dev.

417     Evol. 308:655–668.

418 Mabee P.M., Ashburner M., Cronk Q., Gkoutos G.V., Haendel M., Segerdell E., Mungall C.,

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

420     Trends Ecol. Evol. 22:345–350.

421 Maglia A.M., Leopold J.L., Pugener L.A., Gauch S. 2007. An anatomical ontology for

422     amphibians. Proceedings of the Pacific Symposium on Biocomputing 2007:367-378.

423 Manda P., Balhoff J.P., Lapp H., Mabee P., Vision T.J. 2015. Using the phenoscape

424     knowledgebase to relate genetic perturbations to phenotypic evolution. Genesis. 53:561–

425     571.

426 Manda P., Balhoff J.P., Vision T.J. 2016a. Measuring the importance of annotation granularity to

427     the detection of semantic similarity between phenotype profiles. Proceedings of the

428     International Conference on Biological Ontology 2016.

429 Manda P., Mungall C.J., Balhoff J.P., Lapp H., Vision T.J. 2016b. Investigating the importance

430     of anatomical homology for cross-species phenotype comparisons using semantic similarity.

431     Proceedings of the Pacific Symposium on Biocomputing 2016.:132–143.

432 Midford P.E., Dececchi T.A., Balhoff J.P., Dahdul W.M., Ibrahim N., Lapp H., Lundberg J.G.,

433     Mabee P.M., Sereno P.C., Westerfield M., Vision T.J., Blackburn D.C. 2013. The

434     Vertebrate Taxonomy Ontology: a framework for reasoning across model organism and

435     species phenotypes. J. Biomed. Semantics. 4:34.

436 Mullins P.L., Kawada R., Balhoff J.P., Deans A.R. 2012. A revision of Evaniscus (Hymenoptera,

21

437        Evaniidae) using ontology-based semantic phenotype annotation. Zookeys:1–38.

438    Mungall C.J., Torniai C., Gkoutos G.V., Lewis S.E., Haendel M.A. 2012. Uberon, an integrative

439        multi-species anatomy ontology. Genome Biol. 13:R5.

440    Rámirez M.J., Coddington J.A., Maddison W.P., Midford P.E., Prendini L., Miller J., Griswold

441        C.E., Hormiga G., Sierwald P., Scharff N., Benjamin S.P., Wheeler W.C. 2007. Linking of

442        digital images to phylogenetic data matrices using a morphological ontology. Syst. Biol.

443        56:283–294.

444    Russell J.J., Theriot J.A., Sood P., Marshall W.F., Landweber L.F., Fritz-Laylin L., Polka J.K.,

445        Oliferenko S., Gerbich T., Gladfelter A., Umen J., Bezanilla M., Lancaster M.A., He S.,

446        Gibson M.C., Goldstein B., Tanaka E.M., Hu C.-K., Brunet A. 2017. Non-model model

447        organisms. BMC Biol. 15:55.

448    Ruzicka L., Bradford Y.M., Frazer K., Howe D.G., Paddock H., Ramachandran S., Singer A.,

449        Toro S., Van Slyke C.E., Eagle A.E., Fashena D., Kalita P., Knight J., Mani P., Martin R.,

450        Moxon S.A.T., Pich C., Schaper K., Shao X., Westerfield M. 2015. ZFIN, The zebrafish

451        model organism database: Updates and new directions. Genesis. 53:498–509.

452    Sprague J., Doerry E., Douglas S., Westerfield M. 2001. The Zebrafish Information Network

453        (ZFIN): a resource for genetic, genomic and developmental research. Nucleic Acids Res.

454        29:87–90.

455    The Gene Ontology Consortium. 2017. Expansion of the Gene Ontology knowledgebase and

456        resources. Nucleic Acids Res. 45:D331–D338.

457    Thessen A.E., Bunker D.E., Buttigieg P.L., Cooper L.D., Dahdul W.M., Domisch S., Franz

22

458　　　N.M., Jaiswal P., Lawrence-Dill C.J., Midford P.E., Mungall C.J., Ramírez M.J., Specht

459　　　C.D., Vogt L., Vos R.A., Walls R.L., White J.W., Zhang G., Deans A.R., Huala E., Lewis

460　　　S.E., Mabee P.M. 2015. Emerging semantics to link phenotype and environment. PeerJ.

461　　　3:e1470.

462　Vos R.A., Balhoff J.P., Caravas J.A., Holder M.T., Lapp H., Maddison W.P., Midford P.E.,

463　　　Priyam A., Sukumaran J., Xia X., Stoltzfus A. 2012. NeXML: rich, extensible, and

464　　　verifiable representation of comparative data and metadata. Syst. Biol. 61:675–689.

465　Washington N.L., Haendel M.A., Mungall C.J., Ashburner M., Westerfield M., Lewis S.E. 2009.

466　　　Linking human diseases to animal models using ontology-based phenotype annotation.

467　　　PLoS Biol. 7(11):e1000247.

468　Wiley E.O., Johnson G.D. 2010. A Teleost Classification Based on Monophyletic Groups. In:

469　　　Nelson J.S., Schultze H.P., Wilson M.V.H., editors. Origin and Phylogenetic

470　　　Interrelationships of Teleosts. München, Germany: Verlag Dr. Friedrich Pfeil. p. 123–182.

471　Yoder M.J., Mikó I., Seltmann K.C., Bertone M.A., Deans A.R. 2010. A gross anatomy ontology

472　　　for hymenoptera. PLoS One. 5(12):e15991.