

1 Running Head: EXTRACTING PRESENCE/ABSENCE PHENOTYPES

2

3 Title: Toward Synthesizing Our Knowledge of Morphology: Using Ontologies and
4 Machine Reasoning to Extract Presence/Absence Evolutionary Phenotypes
5 across Studies

6

7 T. Alexander Dececchi¹, James P. Balhoff², Hilmar Lapp², Paula M. Mabee^{1*}

8 *1-Department of Biology, University of South Dakota, Vermillion, SD 57069,*
9 *USA*

10

11 *2-National Evolutionary Synthesis Center, Durham, NC 27705 USA; University of*
12 *North Carolina, Chapel Hill, NC 27599 USA*

13 *Corresponding author

14 Email: alex.dececchi@usd.edu, pmabee@usd.edu

15 ABSTRACT

16 The reality of larger and larger molecular databases and the need to integrate
17 data scalably have presented a major challenge for the use of phenotypic data.
18 Morphology is currently primarily described in discrete publications, entrenched
19 in non-computer readable text, and requires enormous investments of time and
20 resources to integrate across large numbers of taxa and studies. Here we
21 present a new methodology, using ontology-based reasoning systems working

1 with the Phenoscape Knowledgebase (KB), to automatically integrate large
2 amounts of evolutionary character state descriptions into a synthetic character
3 matrix of neomorphic (presence/absence) data. Using the KB, which includes
4 more than 55 studies of sarcopterygian taxa, we generated a synthetic
5 supermatrix of 1051 variable characters scored for 639 taxa resulting in over
6 145,000 populated cells. Of these characters, over 76% were made variable
7 through the addition of inferred presence/absence states derived by machine
8 reasoning over the formal semantics of the source ontologies. Inferred data
9 reduced the missing data in the variable character-subset from 98.5% to
10 78.2%. Machine reasoning also enables the isolation of conflicts in the data, i.e.,
11 cells where both presence and absence are indicated; reports regarding
12 conflicting data provenance can be generated automatically. Further, reasoning
13 enables quantification and new visualizations of the data, here for example,
14 allowing identification of character space that has been undersampled across
15 the fin to limb transition. The approach and methods demonstrated here to
16 compute synthetic presence/absence supermatrices are applicable to any
17 taxonomic and phenotypic slice across the tree of life, providing the data are
18 semantically annotated. Because such data can also be linked to model organism
19 genetics through computational scoring of phenotypic similarity, they open a
20 rich set of future research questions into phenotype to genome relationships.

21

22 Keywords: Ontology, Supermatrix, Inference, Evolutionary mapping,
23 Morphological character, Phenotype, Missing data, Character conflict

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15

The analysis of phenotypic traits in a phylogenetic framework is key to addressing the evolutionary questions posed by an increasingly diverse set of domains. For example, understanding the evolution of pharyngeal jaw mechanics in fishes (Price et al. 2010), identifying phenotype associated genes and regulators in forward genomics approaches (Hiller et al. 2012), exploring the key factors in land plant evolution (Rudall et al. 2013), or discovering the role of phenotypic traits in colonization ability (Van Bocxlaer et al. 2010), all rely on the mapping of phenotypic data to phylogeny. Although robust molecular phylogenies have become easier to generate, more broadly available, and increasingly comprehensive, the phenotypic data on which these studies rely have not.

16
17
18
19
20
21
22

Unlike molecular data, phenotypic data are notoriously time-consuming and complex to generate (Burleigh et al. 2013). Moreover, they are described in a highly detailed free-text format in a distributed literature and have not been available in a computable format (Deans et al. 2012). Researchers seeking to aggregate even the seemingly simple information about the presence and absence of phenotypes across a set of species are faced with a substantial manual extraction and abstraction task (e.g., Stewart et al. 2014). Although

1 assertions of the presence and absence of phenotypes abound in the literature,
2 so do descriptions of the variation in other qualities such as shape, size,
3 position, color, etc. In the case of these qualities, presence and absence must
4 be inferred; from the description ‘posterior flap of adipose fin, free from back
5 and caudal fin’ (Lundberg 1992), the adipose fin would be assumed present.
6 Such detailed data, originally collected for phylogenetic reconstruction or
7 taxonomic identification, are desirable for re-use at the more general level of
8 presence and absence where they pertain to broader questions concerning, e.g.,
9 homoplasy, rates, and correlations of phenotype with environment geography,
10 and genes.

11 Here we show that the presence and absence of phenotypes can be
12 extracted automatically from published detailed phenotype descriptions that are
13 annotated using ontologies. For example, if an author asserts that a fin ray is
14 branched in a particular species, we can use the logic inherent in the
15 corresponding ontology-based expression to infer that the fin ray is present.
16 The power of inference across ontology-based phenotypes (Balhoff et al. 2010;
17 Dahdul et al. 2010a; Mabee et al. 2012) from multiple species and multiple
18 studies enables a substantial reduction in the proportion of missing data in a
19 matrix. We here demonstrate that the logical inferences enabled by ontologies
20 significantly expand the coverage of the data, revealing gaps in phenotype and
21 taxon sampling and revealing data conflict across studies. The methods
22 described here not only allow aggregation of phenotypic data into synthetic
23 supermatrices, but also show the need to more broadly adopt the use of

1 ontology annotation in the morphological literature to facilitate linking and
2 integration with other data such as genetic and developmental data from model
3 organisms (Mabee et al., 2012).

4 METHODS

5 The Phenoscape Knowledgebase (KB) contains ontology-annotated
6 phenotype data derived from published character state matrices from
7 phylogenetic treatments (Fig. 1). Annotations are ontological expressions
8 composed according to the Entity–Quality (EQ) formalism (Mungall et al. 2007,
9 2010), using the Phenex software (Balhoff et al. 2010) as described previously
10 (Dahdul et al. 2010a) (Fig. 1). Anatomical entities are represented by terms
11 from the comprehensive Uberon anatomy ontology for metazoan animals
12 (Mungall et al. 2012; Haendel et al. 2014), which is composed in part from
13 independently developed vertebrate multispecies ontologies (Dahdul et al.
14 2010b; 2012). Phenotypic qualities (presence/absence, size, shape,
15 composition, color, etc.) are drawn from the Phenotype and Trait Ontology
16 (PATO) (Gkoutos et al. 2005). Terms for vertebrate taxa are taken from the
17 Vertebrate Taxonomy Ontology (VTO) (Midford et al. 2013). Every
18 morphological matrix annotated in this way is associated with a single
19 publication in the KB.

20 The KB (Fig. 1) contains a total of 19,024 morphological character states
21 corresponding to 651,660 EQ phenotype annotations for 4,399 extant and
22 fossil vertebrates from 139 comparative studies. It is particularly enriched in

1 the comparative skeletal anatomy for fins, limbs, and their support structures
2 (girdles) of sarcopterygian vertebrates (Supplementary Materials, Table 1), the
3 clade in which the ‘fin to limb’ transition occurred (see Shubin et al. 2014 for a
4 recent discussion). Sarcopterygii comprise slightly greater than half of all
5 vertebrates (Barnosky et al. 2011) and include lobe-finned fishes such as
6 lungfish and coelacanths, and tetrapods including amphibians and amniotes.
7 These richly annotated taxa and phenotypes served as the source data for this
8 investigation.

9 To automate synthesis of supermatrices from the phenotype-by-taxon
10 knowledge in the KB, we created the OntoTrace tool (Fig. 1). OntoTrace
11 accepts as input (1) the targeted anatomical elements (or regions) in the form
12 of a pertinent ontology class or expression (specifically, an OWL class
13 expression), and (2) the taxonomic group(s) (also in the form of an OWL class
14 expression) for which a supermatrix is to be synthesized. OntoTrace first
15 generates a matrix column, and thus a character, for each anatomy ontology
16 class subsumed by the input class expression. Then, for each anatomical
17 character so generated, OntoTrace queries the KB for character states whose
18 EQ annotations logically entail the presence or absence of the respective
19 anatomical element, given the subclass, paratomy, developmental origin, and
20 other axioms provided by the requisite ontologies (see below). The taxa that
21 are associated with those character states and that fall within the input
22 taxonomic group (i.e., are subsumed by the input class expression designating
23 the taxa of interest) are then added to the matrix as rows, and they are given a

1 state of present, absent, or both (as a polymorphism) for the character, as
2 entailed by their respective character states (more precisely, by the EQ
3 annotations for those states). To document the provenance of each synthetic
4 state, all combinations of taxon and published character state supporting the
5 synthetic state value(s), along with references to the respective source
6 matrices (and thus publications), are recorded as metadata for each cell in the
7 synthetic matrix. In addition, OntoTrace determines whether any of the
8 published states supporting a synthetic state directly assert, in the form of
9 their EQ annotations, the presence or absence of the anatomical element, or
10 whether the synthetic state is solely based on logical inference from the
11 supporting states' EQ annotations. Direct assertion of presence/absence here
12 means that the curated EQ annotation(s) for the respective state use the
13 respective character's anatomical structure as the entity ('E'), and one of the
14 terms 'present' (PATO:0000467) or 'absent' (PATO:0000462) as the quality
15 ('Q'). OntoTrace outputs the generated matrix and all metadata in a single file in
16 the NeXML format (Vos et al. 2012) (Fig. 1). OntoTrace is implemented in the
17 Scala programming language, and its source code is freely available under the
18 MIT license on GitHub at <https://github.com/phenoscape/ontotrace>. The source
19 code also contains several ancillary reporting scripts we used to review
20 properties of the matrix (described below). The version of OntoTrace described
21 here has been archived at <http://dx.doi.org/10.5072/zenodo.12705>.

22 To allow manual review of the provenance of the cells in the generated
23 synthetic presence/absence supermatrices, we developed a new interface panel

1 for Phenex, the EQ annotation tool (Balhoff et al. 2010) (Figs.1, 2). Upon
2 selection of a matrix cell in Phenex, the new Supporting State Sources panel
3 displays the list of originally published character states that support the
4 presence/absence state value assignment(s) for the respective taxon, and
5 Phenex highlights in bold those that are considered supporting by direct
6 assertion rather than by inference.

7 To illustrate the properties and value of synthetic morphological
8 supermatrices, we aimed to generate a synthetic presence/absence supermatrix
9 of any anatomical elements that are part of the paired limb, paired fin, and/or
10 the girdle skeletons. We also included elements that are connected to these
11 structures for any sarcopterygian taxa, such as the sternum. To achieve this,
12 we selected anatomical structures using the following OWL class expression,
13 shown below with term labels rather than identifiers for readability:

14 *part_of* some ('paired limb/fin' or 'girdle skeleton') or
15 *connected_to* some ('paired limb/fin' or 'girdle skeleton')

16 The properties *part_of* (BFO:0000050) and *connected_to* (RO:0002170)
17 are from the Relations Ontology (Smith et al. 2005), and the classes 'paired
18 limb/fin' (UBERON:0004708) and 'girdle skeleton' (UBERON:0010719) are from
19 Uberon (Mungall et al. 2012; Haendel et al. 2014). The taxa were selected using
20 the VTO (Midford et al. 2013) term 'Sarcopterygii' (VTO:001464) as input,
21 which permitted us to generate data for taxa annotated to Sarcopterygii or any
22 of its subclasses in VTO. We ran OntoTrace on a Linux-based compute node,

1 using 60 GB RAM, within Duke University's shared high-performance computing
2 cluster, with a build of the Phenoscope KB generated on June 23, 2014.

3 *Entailment of Presence and Absence*

4 The ontologies from which we draw our terms provide a rich context with
5 a community-vetted set of definitions and relationships (structural,
6 developmental) for each entity. The semantics of the OWL ontology language
7 used by Phenoscope, Uberon, and the major model organism ontology
8 communities permits a rich set of inferences to be derived from EQ annotations
9 either in a developmental phenotypic context (as used by model organism
10 databases) or, as seen here (Fig. 3), in an evolutionary phenotypic context. For
11 example a simple EQ annotation may assert that a character state describes an
12 entity 'humerus' bearing a quality 'L-shaped'. A state assignment to a taxon
13 implies that the taxon has a member organism that exhibits a phenotype, that
14 is, has an instance of the class 'L-shaped' that inheres in an instance of the
15 class 'humerus'. Based on this assertion, we can trivially infer that a humerus
16 must be present in the organism. Using an OWL reasoner and additional axioms
17 provided by the Phenoscope KB, more indirect inferences of presence or
18 absence can be made as well, which essentially result from the anatomical
19 knowledge expressed within the Uberon ontology (Balhoff et al. 2014).

20 *Presence.*—To query character states that denote presence of a given
21 structure, OntoTrace retrieves phenotypes from the KB that are subsumed by
22 the expression '*implies_presence_of* some <entity>'. *Implies_presence_of* is an

1 OWL property that unifies the various means by which the presence of a
2 structure can be inferred (see Balhoff et al. 2014 for details). For example, a
3 quality that inheres in a structure *implies_presence_of* that structure (Fig. 3).
4 Presence is also inferred for any structures of which that structure is a part or
5 from which it develops. The presence of a ‘humerus’ implies the presence of a
6 ‘forelimb’ and a ‘forelimb skeleton’, of which Uberon asserts it to be a part. The
7 presence of a ‘forelimb’ also implies the presence of a ‘forelimb bud’, because
8 Uberon asserts that the former develops from the latter (Fig. 3).

9 *Absence.*—To query character states that denote absence of a given
10 structure, OntoTrace retrieves from the KB those phenotypes that are
11 subsumed by the expression ‘lacks_all_parts_of_type and *inheres_in* some
12 multicellular_organism and *towards* value <entity>’. Similar to ‘presence’, the
13 KB makes use of chains of ontological relationships to infer which other
14 structures must be absent as the consequence of the absence of a given
15 structure (Fig. 3) (see Balhoff et al. 2014 for details). For example, the absence
16 of a ‘forelimb’ entails the absence of a ‘humerus’.

17 *Identifying Conflicts*

18 When a taxon is inferred to exhibit both presence and absence for a
19 particular structure, it indicates either a polymorphic condition in the taxon, or
20 the fact that the supporting original character states, or more precisely, the EQ
21 annotations made for them, conflict with each other. Polymorphic synthetic
22 state values were considered as reflecting actual polymorphism, and thus

1 excluded from further review, if both presence and absence are directly
2 asserted by supporting character states associated with a single source matrix
3 (and thus the same publication). To aid manual review of the remaining
4 conflicts, we created a script that reported for each conflict the taxon, the
5 entity that was polymorphic (i.e., had conflicting values), and whether the
6 presence and absence values were supported by direct assertion or inference.
7 This reporting script can be found in the OntoTrace source code repository.
8 Provenance of conflicting data can be viewed in Phenex (Fig. 2).

9 *Identifying Isomorphic Synthetic Characters*

10 To examine possible dependence across characters in the synthetic matrix
11 as a consequence of assertions in the ontologies, we used a script to report
12 each cluster of characters (i.e., anatomical entities) that were identical in their
13 taxonomic distribution of values. In other words, all identical character columns
14 were collected into clusters; we term these clusters ‘isomorphic characters’.
15 Further, for each of the anatomical entities comprising each cluster, the script
16 reported whether the matrix contained any direct presence/absence assertions
17 for that character, or if it was included in the matrix solely through inference.
18 Code for this report can be found in the OntoTrace source repository. To aid in
19 characterization of the ontological dependence of isomorphic characters, we
20 used a script to generate an ontology of ‘presence classes’: For each anatomical
21 entity ‘X’ from the Uberon ontology, we generated a corresponding class with
22 the logical definition *‘implies_presence_of some X’*. We classified these
23 expressions using the ELK reasoner (Kazakov et al. 2012, 2013) within the

1 Protégé ontology editor, and used the Protégé DL Query panel to check for
2 inferred equivalency between presence expressions.

3 *Other Reporting Queries*

4 The number of published character states that entail the presence or
5 absence for selected sets of entities and taxa was reported using queries to the
6 Phenoscape KB implemented as a script included within the OntoTrace source
7 code. Specifically, for each taxon and entity, we queried for states that were
8 annotated with phenotypes that entailed either the presence or absence of the
9 entity.

10 Using a SPARQL query to the Phenoscape KB, we counted the number of
11 published matrices in which each sarcopterygian taxon in the KB is included. An
12 additional SPARQL query was used to report, for each published matrix in the
13 KB, the number of taxa, characters, states, and phenotypes associated with
14 annotations relevant to structures of the fin or limb. These queries are included
15 within the OntoTrace code repository.

16 RESULTS

17 OntoTrace aggregated, as described above, the KB's morphological
18 phenotype data on paired limb, paired fin, and/or the girdle skeletons for all
19 sarcopterygian taxa into an entity-by-taxon matrix of 1,759 synthetic
20 presence/absence characters by 1,052 taxa, in the form of an XML file in
21 NeXML format (available from Dryad). The 55 studies from which data were

1 synthesized in this manner are summarized in Supplementary Materials Table 1.
2 The data from these papers that relate to fin, limb, girdle and their parts total
3 2,588 text-based character states from 1,195 individual published characters,
4 scorable for 1,052 sarcopterygian taxa.

5 Out of the 1,759 generated synthetic characters, 639 were variable, i.e.,
6 included both presence and absence states. Of these, 488 characters (76%)
7 were variable only due to the use of inference: 442 variable characters were
8 composed of inferred data alone; 12 were made variable by inferred absence,
9 and 34 by inferred presence. In the matrix subset comprising the variable
10 characters there are 146,451 populated cells, which constitute 21.8% of all
11 cells. Directly asserted data accounted for only 9,948 (6.8%) of the populated
12 cells, or 1.5 % of all cells in the subset; in contrast, inferred data represent
13 93.2% of the populated cells (Fig. 4). Of the 1,051 taxa in the subset, 13%
14 (136 taxa, see Supplementary Materials, Table 2) are included in the matrix
15 solely on the basis of inferred data. Taxa for which the source matrices contain
16 no direct assertions about presence or absence of any fin/limb entity can
17 nonetheless be included in the synthetic presence/absence matrix if they have
18 EQ phenotype annotations that imply presence or absence of a fin/limb entity.
19 For example the theropod dinosaur taxon *Sinosauropteryx prima* in our data is
20 derived from a single source (Sereno et al. 2009), where it was not scored for
21 any presence/absence characters. Its inclusion in the synthetic supermatrix
22 comes solely from character states such as ‘increased scapular blade width’ and
23 ‘poorly differentiated humeral head form’, because these imply the presence of

1 a scapula and humerus, respectively.

2 After excluding polymorphisms directly asserted within a single source
3 (see 'Identifying conflicts' above), we identified 774 cells (of the 146,451
4 populated ones) as stating both presence and absence (0/1) of a character, for
5 99 synthetic characters and 297 taxa. These included 135 conflicts between
6 direct assertions (that were made in different source publications), 565
7 conflicts between direct assertions and inferred states, and 74 conflicts
8 between inferred states (Supplementary Materials, Table 3).

9 We also identified 93 clusters of characters in the synthetic supermatrix
10 that were isomorphic, i.e., identical in their distribution across taxa and to one
11 another, but variable. These correspond to 85,813 cells in the synthetic
12 supermatrix (Supplementary Materials, Table 4), almost 59% of the (populated)
13 cells. To better characterize these clusters as to their ontological basis, we
14 examined which of them fall into equivalence chains of implied presence and
15 absence. More specifically, two synthetic characters with anatomical entities X
16 and Y, respectively, for which a reasoner infers equivalence between the logical
17 definitions '*implies_presence_of* some X' and '*implies_presence_of* some Y' will
18 necessarily be found isomorphic in their distribution of presence and absence.
19 For example, 'presence of pedal digit 2' is inferred as equivalent to 'presence of
20 pedal digit 2 digitopodial skeleton'. Twenty-one of the 93 clusters (23%) were
21 of this kind. Another 6 (6%) were found to be clusters of anatomical parts and
22 the entities that contain them. Clusters can also arise from co-asserted entities
23 (e.g., when a single character state includes multiple entities, such as 'pedal

1 digits 6, 7, and 8 present', which will result in 3 EQ annotations, one for each of
2 the 3 digits). There were 3 such clusters (3%). Most of the clusters were
3 composed of inferred data only. Of these, 63 (68%) resulted from chains of
4 inference from multiple entities that were different for each cluster.

5 The number of source matrices from which a particular taxon was
6 sampled ranged from 1 to 16 (Supplementary Materials, Table 5). 813 (77.4%)
7 of the sampled taxa are at the species rank, with the remainder distributed
8 across higher-level ranks (Supplementary Materials, Table 5).

9 The number of published character states that entail the presence or
10 absence for selected parts of the fin and limb was used to generate a figure
11 showing their distribution across taxa along the fin to limb transition (Fig. 5).

12 DISCUSSION

13 The first step in scaling up the exploration of phenotypic patterns in an
14 evolutionary context is to render phenotypic descriptions of species in a form
15 amenable to large-scale computational integration, linking, and mining. How this
16 is possible has recently been shown in a series of papers from the Phenoscope
17 group (Dahdul et al. 2010a; Mabee et al. 2012). Here we demonstrate that,
18 using computable phenotypes from a datastore representing the cumulative
19 effort of experts across a broad taxonomic scale, synthetic supermatrices for
20 presence/absence phenotypes can be automatically assembled for user-
21 designated slices of the taxonomic and anatomical corpus.

1 Bringing together phenotypic data from across multiple studies manually,
2 and synthesizing them in a form amenable to computational analysis, is a non-
3 trivial exercise. Manual concatenation of phylogenetic matrices, for example,
4 necessarily involves the time-consuming process of identifying and eliminating
5 character redundancy (e.g., Gatesy et al. 2002; Gatesy and Springer 2004; Hill
6 2005). Ascertaining the presence or absence of a morphological feature
7 requires an additional effort to reason from text that may only incidentally
8 describe an aspect of it. As a consequence, the ability of scientists to hand-
9 assemble data across studies is severely hampered by the difficulty to compute
10 on taxa and morphological data. Our work shows that computable phenotypes
11 not only enable automatic consolidation of character states into non-redundant
12 presence/absence assertions, but they enable inference of presence/absence
13 generalizations. Our method makes data reuse by non-experts not only more
14 efficient, but also reduces the risk for error and expands the phenotypic and
15 taxonomic coverage of the original data. In so doing, it can open new
16 possibilities for data analysis, in particular if phenotypes are linked with genes
17 and other data through their shared ontological context (Mabee et al. 2012).

18 *Inference Expands Data: Filling in the 'Unknown Knowns'*

19 Our results demonstrate that inference can play a profound role in
20 supplementing the taxonomically sparse phenotype assertions across taxa (Fig.
21 4). We found that 76% of the variable characters in the synthetic supermatrix
22 were made that way through inference, meaning that at least one of their two
23 states (presence or absence) is based solely on inferred data. For an individual

1 feature such as the skeleton of digitopodium (the collection of skeletal
2 elements encompassing the digits, i.e., the metacarpals/tarsals and phalanges),
3 the number of inferred assertions (7751 annotations for 718 taxa) is 38 times
4 higher and spread over 7 times more taxa than direct assertions (201
5 annotations for 103 taxa). An individual taxon can have multiple sources of
6 inference for an individual entity, depending on the number and nature of the
7 annotated characters that relate to that taxon and entity. For example,
8 *Acanthostega* has 5 directly asserted and 30 inferred sources with character
9 states that entail presence or absence of ‘skeleton of digitopodium’.

10 At the most basic level, aggregation of and inference on phenotype data
11 allows users to supplement large amounts of missing data computationally,
12 without extensive manual literature research. As Hiller et al. (2012) show,
13 simply knowing in which taxa a phenotypic trait is present or absent across a
14 taxonomic range in which the trait underwent evolutionary change can enable
15 entirely new insights into the developmental genetics of the trait. While the
16 overall goal of our method, filling in data that is not directly asserted, is similar
17 to imputation, our method differs substantially from this technique. Regression-
18 based imputation practices for finding the ‘invisible fraction’ (Grafen 1988) use
19 a probabilistic models to reconstruct the ‘unknown knowns’, while our method
20 bases its reconstructions on predefined logical axioms in the ontology.
21 Imputation methods can be effective at reducing gaps in quantitative data sets
22 (Nakagawa and Freckleton 2008, 2011; Swenson 2014), however they are not
23 applicable to qualitative matrix data. Our use of inference to extract unstated

1 knowledge about the presence and absence of traits allows reconstruction of
2 missing values without resorting to statistical parameters that may change
3 across phylogeny. Though we restrict ourselves to a simple set of relational
4 rules for entities and their parts that are uniform across metazoans (i.e., the
5 humerus, when present, is always *part_of* the forelimb), the results of the
6 logical reasoning methods used here are very powerful.

7 The supermatrix technique is a total evidence approach in systematics
8 (Kluge 1989), where different data sets and types are combined into a single
9 ‘supermatrix’ of unique taxa (Sanderson 1998; de Queiroz and Gatesy 2007).
10 Inevitably, component data sets overlap, but incompletely, resulting in many
11 taxa lacking data for many characters. In the realm of molecular data matrices,
12 there have been two approaches to deal with the missing data problem: (1)
13 leave all taxa separate and code the unavailable characters as missing; or (2)
14 reduce missing data by making composite taxa at a level for which monophyly is
15 assumed *a priori*. The former of these may lead to loss of resolution but not
16 necessarily misleading relationships (Kearney 2002; Kearney and Clark 2003;
17 Wiens and Morrill 2011; Wiens and Tiu 2012). The latter, composite taxa, may
18 lead to misleading phylogenetic results (Malia et al. 2003). Whether
19 supermatrices are sequence or morphology based, they typically involve a lot of
20 missing data. Molecular supermatrices may include over 70% missing data (de
21 Queiroz and Gatesy 2007; Fabre et al. 2009; Hejzol et al. 2009; Hedtke et al.
22 2013), and a morphological supermatrix assembled by Ramírez et al. (2007)
23 had 94% missing data. By comparison, missing data in the variable character-

1 subset of the synthetic supermatrix we created amounted to 98.5% without
2 applying inference, and applying logical inference reduced this fraction to
3 78.2%.

4 *Characteristics and Application of Isomorphic Synthetic Characters*

5 A considerable fraction (nearly 60%) of the populated cells in the variable
6 subset of the synthetic supermatrix are characters that are part of one of 93
7 isomorphic clusters. That is, they corresponded to entities whose
8 presence/absence distributions are identical across taxa (Supplementary
9 Materials, Table 4). That our matrix synthesis method generates isomorphic
10 clusters is expected, because presence/absence reasoning uses axioms in the
11 Uberon anatomy ontology about parthood and developmental precursor
12 relationships (Balhoff et al. 2014). This will necessarily result in clusters
13 composed of an asserted entity and its containing (for presence) or contained
14 (for absence) classes, and/or developmental precursors (for presence) or
15 derivatives (for absence). For example, for a character of 'femur bone', a state
16 value of present will induce the same state value for characters 'hindlimb',
17 'limb', 'femur cartilage element', 'femur pre-cartilage condensation', and so on
18 (see Fig. 3). We found that 10% (9 of 93) of the isomorphic clusters were of
19 this kind. Most clusters were composed of only inferred data, indicating that the
20 underlying original character states did not directly assert presence or absence.
21 Although some of these (21 clusters) could be identified as the consequence of
22 logic equivalence chains for implied presence or absence (see Results), the
23 majority (63 clusters) resulted from various chains of inference from multiple

1 entities with no obvious repeating patterns. For example, the taxonomic
2 distributions of presence for 'nail', 'dorsal skin of digit', 'distal limb
3 integumentary appendage' and 'digit skin' are identical, even though the
4 ontological presences of these entities are not inferred to be equivalent.

5 Whether and what value or impact these isomorphic clusters have will
6 depend on the goals of the researcher using these data. Perhaps the broadest
7 and most forward-looking applications for phenomic data involve understanding
8 trait evolution in relation to other phenotypic traits, environmental factors, and
9 aspects of development and related genomic features. For research questions
10 such as these, the biological knowledge revealed by a cluster of correlated
11 characters may have substantial value in supplementing the input data. For
12 example, the presence of developmental precursors (femur cartilage, femur
13 condensations) entailed by an entity's presence (e.g., the femur bone) may
14 involve different genes and networks relevant to a developmental biologist. One
15 can also envision research questions where the inferred data hold value that the
16 asserted data do not. For example, the inferred presence and absence of
17 'hindlimbs' would be valuable for studies examining correlations of habitat and
18 locomotion, whereas the phenotype assertions that entailed them may not be
19 directly related to locomotion (e.g., toenail color). Further, the knowledge
20 structure that is laid out in sets of isomorphic characters may be of benefit as
21 approaches to computationally dissecting out the expression of genes and their
22 regulators (Hiller et al. 2012) are scaled up.

23 For researchers interested in using these matrices for phylogenetic

1 reconstruction, caution must be exercised. The character dependency implied
2 by a significant fraction of isomorphic (even if otherwise variable) characters
3 suggests that synthetic matrices, at least in the form of presence/absence-only
4 data, are not immediately suitable for phylogenetic reconstruction. As discussed
5 above, observed isomorphism does not necessarily imply logical equivalence,
6 and hence whether characters should be merged or not due to putative
7 dependency would need to be carefully examined for each case. For example,
8 'nail' and 'dorsal skin of digit', though isomorphic in their taxonomic distribution
9 in this dataset, have different developmental bases, and can thus be argued to
10 not be dependent.

11 More generally, the degree to which phylogeny can be recovered from
12 binary presence/absence data alone has, to our knowledge, not been
13 investigated. Certainly presence/absence data are common in morphological
14 datasets; Sereno (2009) gives a figure of 25%. However the phylogenetic
15 resolution attained in these studies require variation in other qualities (size,
16 shape, texture, color, etc.). The ontological methods used here reduce data
17 from these other qualities to presence/absence, thus changing the phylogenetic
18 level at which the information is relevant. For example, if variation in vertebral
19 shape across a set of taxa is reduced to the inference that vertebrae are
20 present in these same taxa, it no longer contains information to resolve them.
21 However the presence of vertebrae is informative for resolving taxa at a higher
22 level (i.e., as members of the monophyletic clade, Vertebrata). Though this
23 issue will require further examination, it is likely that the inferred

1 presence/absence data will only support the monophyly of more inclusive clades
2 than the original assertions.

3 *Distribution of Data Across Taxa and Anatomical Regions*

4 The paired appendages are ostensibly one of the most intensely studied
5 aspects of anatomy in vertebrates, and yet quantifying the data available for
6 them has not previously been possible. The methods presented here readily
7 enable this, including visualizing how our knowledge of morphology, whether
8 expressly stated or implied, is distributed over taxonomic and anatomical space
9 (Fig. 5). This can then be used to pinpoint the taxonomic groups and the parts
10 of the anatomy for which data are sparse or lacking, allowing potential reasons
11 and remedies to be considered. One should note in this context that availability
12 and lack of data for an anatomical feature in a taxonomic group should not be
13 expected to coincide with presence and absence, respectively, of the feature in
14 said group. Figure 5 illustrates this, for example, for digits in the lungfish
15 *Dipterus*. Lungfishes do not have digits, yet due to assertions about their
16 absence in this taxon (Zhu et al. 1999; Swartz 2012) data about digits in
17 lungfishes are available.

18 For the matrix we synthesized for the evolution of vertebrate fin/limb
19 morphology, the gaps in the data may be primarily attributable to the following
20 two factors. One, most taxa studied in the fin to limb transition are fossils and
21 thus restricted to a few, often partial, specimens. These taxa may also be
22 unscorable for certain entities due to primitive absence (e.g., the ilium, ischium

1 and pubis of the pelvis are not present in basal non-tetrapod taxa). Two, the
2 taxa and anatomical elements used for study are unequally sampled. As is
3 evident in Figure 5, there are much less data about the hindlimb relative to the
4 forelimb in basal tetrapods, which cannot be explained by hindlimb specimens
5 being unavailable or far less preserved in the fossil record for the respective
6 taxa than their forelimbs. Hence, other explanations are needed. Perhaps the
7 differences could be due to more variability, and therefore more character data,
8 in the forelimb than the hindlimb during the fin-limb transition, which would be
9 consistent with the “front wheel drive’ hypothesis, which posits that the fin to
10 limb transition was driven primarily by changes in the forelimb (Shubin et al.
11 2014). Alternatively, the difference could be a result of sampling bias caused by
12 the larger size of the ancestral forefin and the interconnectedness of the girdle
13 skeleton with dermal skull elements.

14 Regardless of what is really behind the difference, our results illustrate
15 how the ability to visualize the uneven distribution of knowledge can reveal far
16 more than simply the existence of bias. Gaps in morphological knowledge, such
17 as here the phenotypic evolution of the hindlimb, can present major challenges
18 for understanding the origins and evolution of novel features (Shubin et al.
19 2014), and the ability to synthesize knowledge on a large scale can focus future
20 studies on filling in gaps.

21 *Quantification of Taxon Scoring*

22 As a consequence of the obstacles to integrating morphological character

1 data, it has been nearly impossible to assess quantitatively the differential
2 sampling of taxa and anatomy across studies. This, too, is readily enabled by
3 the methods described here. As an example, we examined how frequently
4 individual taxa had been scored for fin and limb phenotypes in the generated
5 synthetic supermatrix. Because of the logistic efforts necessarily involved in
6 morphological data collection (specimen preparation, museum collection visits,
7 etc.), the taxonomic sample of species that an investigator can examine is
8 limited, and some taxa are more readily available for study than others. In our
9 dataset 70% of the taxa in the synthetic supermatrix were connected to only a
10 single publication record (Supplementary Materials, Table 5). For taxa having
11 more than one source publication, the proportions drop rapidly: 12% and 7%
12 are found in two and three publications, respectively, and less than 2% of the
13 taxa are scored in seven or more publications. A single taxon, *Acanthostega*, a
14 well-preserved exemplar taxon in the fin to limb transition, holds the maximum
15 number of 16.

16 However, this distribution, and in particular the high proportion of single-
17 source publication taxa, is unlikely to be representative of the vertebrate
18 comparative fin/limb morphology literature as a whole. This is because the
19 publications we chose for phenotype annotation treat mostly non-overlapping
20 sections of the vertebrate phylogeny, and thus a high fraction of taxa with a
21 single publication source is a consequence of our experimental design. If we
22 consider only the data for basal sarcopterygians relating to the fin to limb
23 transition (Supplementary Materials, Table 1), the proportion of taxa with only a

1 single publication source drops to 34%. However, when considering the fraction
2 of taxa whose morphological features have been scored independently only
3 once, this figure is likely an underestimate. Some of the publications in this
4 subset of the supermatrix share co-authors, and many characters are recycled.
5 A more thorough study of independence and depth of evidence across the
6 dataset was beyond the scope of this study, but our results illustrate how our
7 methods would readily enable such an analysis.

8 *Conflicting Data Revealed*

9 When authors reuse characters from previous works, encountering, and
10 resolving coding conflicts is an important part of the process to ensure
11 phylogenetic relationships are as accurate as possible (Harris et al. 2007).
12 Character conflicts are often difficult to spot by hand, yet the protocols authors
13 follow for identifying, adjudicating, and resolving conflicts are rarely reported
14 beyond a high-level summary. The presented supermatrix synthesis approach
15 immediately reveals conflicting phenotypes, here in the form of an anatomical
16 feature having state values of both present and absent (0/1) for the same
17 taxon. We found 774 such cells (0.5%) among the 146,451 populated cells,
18 excluding directly asserted polymorphisms, which we defined as those that
19 trace back to direct assertions of both states in the same source matrix (see
20 Methods). How this level of character conflict compares to what has been
21 observed previously is difficult to assess, because in previous studies in which
22 morphological matrices have been concatenated manually (see O'Leary et al.
23 2013; Sigurdson and Green 2011), the resolution of conflicts is not reported in

1 a quantitative manner. However, in a consensus morphological matrix for
2 turtles, Harris et al. (2007) reported <2% cells with conflict (out of 4872 total
3 cells), which is similar in magnitude to our finding.

4 One of the major advantages of the synthesis approach we present is not
5 only that the extent of character conflict can be quantified quickly, but also
6 that detailed reports about the provenance of all conflicting data can be
7 generated automatically. This greatly aids review, and where possible, resolution
8 of these data by experts. There are a number of different causes for an
9 observed conflict, only some of which are correctable errors; determining the
10 cause for a given character conflict requires careful examination. A trivial case
11 stems from the fact that many or even most matrices are not yet archived in
12 digital repositories (Stoltzfus et al. 2012; Drew et al. 2013), and errors could
13 be a result of their required manual digitization. These will be reduced by the
14 increasing push for digital archival of matrices upon publication. More
15 substantive conflicts however result from differing author assertions that may
16 stem from observations of different (and differing) specimens or different
17 interpretations of the same material. Additionally, the conceptualization of the
18 character by the original author, and the terminology used for its description,
19 may have consequences beyond the confines of the original state structure
20 when annotated with ontology terms that have logical implications, leading to
21 conflicting results. A discussion of conflicts in relation to their bases in assertion
22 and/or inference follows, with examples from the dataset we generated. It is
23 worth noting that for conflicts due to correctable errors, our fully computational

1 approach to matrix synthesis has the advantage that once the errors are
2 addressed in the KB, the corresponding conflicts are eliminated from any
3 supermatrix subsequently generated from it.

4 *Conflicts between Asserted Character States*

5 The conflicts that are most readily traced to their cause are those
6 between authors who differently assert the presence and absence of an
7 anatomical structure. These comprise a relatively small proportion of the
8 conflicts (17%). Some of these discrepancies arise as new observations are
9 made, e.g., from new specimens that reveal formerly poorly known anatomy.
10 For example, Zhu and colleagues (1999), scored the humerus of *Strepsodus*, a
11 rhizodontid fish, for the presence of distinct supinator and deltoid processes.
12 Based on new fossil material for rhizodontids, the humerus morphology was re-
13 evaluated by Jeffery (2001) who concluded that in *Strepsodus* and other
14 rhizodontids distinct supinator and deltoid processes are absent, thus
15 generating the conflict observed in our dataset.

16 Sometimes, the basis of conflict between original author assertions is not
17 as readily traceable. For example, in the fossil literature it is not uncommon that
18 not all of the specimens examined in relation to each operational taxonomic unit
19 (OTU) are reported. Even when specimens are listed comprehensively, the
20 reasons for conflicts are sometimes difficult or even impossible to deduce from
21 the published literature alone. For example, Ruta et al. (2003) state that
22 accessory foramina (passages for blood vessels) are absent in the humerus of

1 the fossil amphibian *Sauroploera*, but later Ruta (2011) scored these foramina
2 as 'present'. As it does not appear from his documentation that different
3 specimens were examined, this leaves re-examining the specimens or
4 communicating with the authors as the only resort to resolving the conflict.
5 Such differences in scoring are a challenge for both manual and machine
6 concatenation of these data, but they are to be expected, as authors not only
7 have access to different materials over time, but will also sometimes vary in
8 their interpretation of structures. The presented matrix synthesis method
9 cannot reduce or eliminate them, but it is able to readily pinpoint candidates for
10 investigation, including by way of computationally (and thus automatically)
11 generated reports.

12 *Conflicts between Asserted and Inferred Character States*

13 The most frequent conflicts (73%) occur between asserted and inferred
14 data. These are arguably much less obvious from manual analysis than the
15 detection of conflicting assertions. An example comes from a recent large-scale
16 examination of tetrapod limb evolution, focused on the transitional fossil
17 *Tiktaalik roseae*, which is described as having a 'poorly developed' scapula blade
18 (Ruta 2011). This assertion results in its inferred presence in the synthetic
19 supermatrix (Fig. 1). The scapula blade, however, is directly asserted to be
20 absent in *Tiktaalik roseae* by Swartz (2012) and Clack et al. (2012). Regardless
21 of what is at the root of this conflict (different specimens, different
22 interpretations of morphology, polymorphism, etc.), the value of our method is
23 that it makes the discrepancies in the literature evident.

1 *Conflicts between Inferred Character States*

2 The fewest conflicts (10%) are generated between data based on
3 inference alone. For instance in the frog *Bombina variegata*, the ilial
4 protuberance is inferred absent based on the assertion that the ilial shaft is
5 absent (Fabrezi 2006), of which the ilial protuberance is a part. The presence of
6 the ilial protuberance is inferred from two assertions regarding its shape, i.e.,
7 ‘not knobbed distally’ and ‘broad and low rounded’ (Cannatella 1985), thus
8 generating a conflict. Identifying the condition(s) in this species is beyond the
9 scope of this paper, but it would likely require the user to analyze the
10 supporting specimens from the original sources. Again, the value of our method
11 is that it reveals the conflicts, here from inference alone, which particularly in
12 this case would be difficult to ascertain manually.

13 *Conflicts from Author Character Structure and Scoring*

14 Some ‘false’ conflicts resulted from the idiosyncratic character
15 construction and scoring practices by authors, and also limitations of the KB.
16 For example, a conflict is automatically generated when an author creates a
17 character state that is a disjunction of absence and one or more other qualities
18 that entail presence. For example, the character, ‘ectepicondyle’ with the state:
19 ‘low, indistinct or absent’ (Laurin and Reisz 1997) is intended to reflect the
20 variability present across the taxa. Yet this wording does not allow the reader to
21 differentiate whether this represents polymorphism within species (i.e., different
22 states in different individuals of a single species), or whether the set of species

1 to which the description applies has multiple states (i.e., one state in one
2 species, a different one in another). In this case, because ‘low’ implies the
3 presence of an ectepicondyle, it is automatically shown as in conflict with the
4 same authors assertion of absence. This illustrates how ambiguity in how an
5 author constructs character states can limit or even preclude the utility of their
6 data in other contexts.

7 This particular type of conflict also reflects a problem with our annotation
8 methodology, which does not allow combining phenotypes using ‘or’ as stated
9 by the author. Instead, the system applies all of the described phenotypes to
10 taxa with this state. The resulting conflict could be eliminated within the
11 semantic model by representing the annotation as a logical union of the
12 phenotypes. However, an assertion that a given taxon has an instance of
13 phenotype ‘A or B’ prevents a machine reasoner from applying any inferences
14 based on either A or B to the taxon, because it cannot know which of the
15 classes of phenotypes the taxon actually has. Thus, annotating a state with a
16 ‘union phenotype’ would effectively prevent this character state from
17 contributing to inference of presence/absence, even if such an approach would
18 more accurately reflect the knowledge asserted in the original paper.

19 Another source of error stems from character constructions that involve
20 phenotypes of anatomical elements that are more complex than simple
21 presence/absence, but are applied to taxa to which strictly speaking they don’t
22 apply. For example, the frog *Rhinophrynus dorsalis* is asserted to lack a sternum
23 (Cannatella 1985). Yet in the same study the epicoracoid bone is scored in this

1 taxon as ‘not fused to sternum’, from which a machine reasoner, and arguably
2 also a human reader, would infer that a sternum is present. If an author scored
3 this as ‘not applicable’ in the original matrix, the logical error (inferred presence)
4 would be avoided.

5 Finally, inattention to the semantics of anatomical terminology can lead to
6 incorrect and conflicting assertions. For example, while there is clearly a deep
7 homology across Sarcopterygii between distal fin (‘paired fin radials’) and distal
8 limb (‘digits’) skeletal elements (Johanson et al. 2007), they are generally
9 considered distinct. Yet in the synthetic matrix, some limbed tetrapods are
10 inferred to possess both radials and digits. This inference was generated from
11 several limbed, and potentially terrestrial tetrapod taxa such as *Acanthostega*,
12 *Dendrerpeton*, and *Silvanerpeton*, scored as possessing ‘jointed radials’ (Swartz
13 2012). Thus the presence of radials is inferred for these taxa, while
14 simultaneously digits were directly asserted for them (Ruta 2011). Perhaps
15 Swartz (2012) used ‘radial’ to encompass all acropodial elements because there
16 is simply not a more encompassing anatomical term that applies to these distal
17 skeletal structures across the taxonomic breadth of vertebrates. This use of
18 ‘radial’, however, conflicts with its general usage in the literature (Ruta 2011)
19 as well as genetic data concerning the distinctness of digits (Davis 2013;
20 Woltering et al. 2014). Referencing and applying a standardized vocabulary in
21 character descriptions would resolve this type of conflict (Seltmann et al.
22 2012).

23 As described above, author-generated conflicts pose a problem to the

1 effort of automatic integration of these manually annotated character data.
2 Because they are idiosyncratic and difficult to detect until integration, there is
3 little possibility to create filters that automatically correct for these types of
4 errors. We suggest it is better to work to amend character construction
5 practices, and working toward a future in which characters are constructed in
6 computable form *a priori*, than trying to address them *post hoc*. This will likely
7 become even more important as text markup of phenotypes and other concepts
8 is automated, leaving little margin for human curator correction of
9 inconsistencies.

10 *Improved Annotation and Curation Standards*

11 The annotation practices that guided the Entity–Quality assignments to
12 the character data in Phenoscape were designed to capture the rich anatomical
13 detail and differences among taxa, as described by taxonomic experts.
14 Combining and reasoning across the annotations in this study, however, cast
15 these data in different relief, in some cases revealing conflicts that were the
16 result of inappropriate annotation of author statements. Resolving these, and
17 generalizing the issues where possible, enabled us to improve and expand the
18 anatomy and quality ontologies, the annotations, and the phenotype curation
19 guidelines (http://phenoscape.org/wiki/Guide_to_Character_Annotation). For
20 instance, it appeared that data from a single paper conflicted in whether or not
21 the fish *Onychodus* possessed a postcleithrum (Cloutier and Arratia 2004).
22 Investigation revealed that the authors directly asserted the absence of this
23 bone; an inferred but incorrect presence resulted from a mistake in annotating

1 'presence of a postcleithral scale' as "dermal scale' and *part_of* some
2 postcleithrum, present'. The postcleithral scale, however, is a separate type of
3 scale and not a part of the postcleithrum bone. In this case we added a new
4 entity 'postcleithral scale' to the Uberon ontology as a type of 'scale', the
5 feature was re-annotated, and the conflict thus removed.

6 FUTURE DIRECTIONS

7 The approach and methods demonstrated here to compute synthetic
8 presence/absence supermatrices are applicable to any taxonomic and
9 phenotypic slice across the tree of life, provided these data are semantically
10 annotated. Scaling up annotation to this level, however, will require significant
11 effort, including the development of semiautomated methods for marking up
12 free-text descriptions (e.g., Cui 2012; Arighi et al. 2013); provisioning of
13 community phenotype ontologies to accommodate the diversity of taxa and
14 evolved anatomies and qualities (Gkoutos et al. 2005; Haendel et al. 2014); and
15 faster and more efficient methods for reasoning across these substantially
16 larger data in knowledgebases. Another challenge lies in developing methods to
17 aggregate 'non' presence/absence phenotypes, i.e., those features varying in
18 qualities such as size, shape, color, texture, etc., into a matrix format, which will
19 require sophisticated algorithms for automating consolidation of synthetic
20 character states. Additionally, new methods are required to integrate
21 taxonomically-heterogeneous supermatrix data with user-specified trees.
22 Because the phenotype data are asserted at multiple taxonomic levels (i.e., to
23 species, genera, families, etc.), current methods for their optimization and

1 visualization along a tree are limited.

2 CONCLUSIONS

3 The phenotypic features that characterize and define evolutionary groups are
4 currently scattered across the dispersed literature of comparative biology, often
5 in character-by-taxon matrices for small sets of taxa. The difficult and time-
6 consuming manual aggregation of these data reduces their reuse. Here we
7 demonstrate that when phenotypes are ontology-annotated, their presence and
8 absence can be automatically integrated into synthetic character matrices. We
9 found that inference plays a profound role in supplementing the taxonomically
10 sparse phenotype assertions across taxa, in our case reducing the missing data
11 in the variable character-subset from 98.5% to 78.2%. Moreover, 76% of the
12 variable characters were in fact made variable through the addition of inferred
13 presence/absence states. Equally important, this automated method results in
14 immediate isolation of character conflicts and detailed reports about their
15 provenance. This capability, if available broadly, will greatly aid experts in data
16 review, and where possible, conflict resolution. Finally, machine reasoning
17 enables quantification and new visualizations of the data, as demonstrated here,
18 allowing the identification of character space that is undersampled across the
19 fin to limb transition.

20 FUNDING

21 This work was supported by National Science Foundation collaborative grants
22 (DBI-1062404, DBI-1062542) and the National Science Foundation National

1 Evolutionary Synthesis Center (NESCent) (EF-0905606).

2

3 ACKNOWLEDGMENTS

4 We thank Phenoscope collaborators, including D. Blackburn, W.M. Dahdul, P.
5 Manda, C. J. Mungall, and T.J. Vision for comments and advice that improved
6 this work. We also acknowledge N. Ibrahim and W.M. Dahdul for annotation of
7 some of the data analyzed here, as well as M.A. Haendel and C. Mungall for
8 guidance and help with ontology development.

9

10

11 REFERENCES

12 Arighi C.N., Carterette B., Cohen K.B., Krallinger M., Wilbur W.J., Fey P., Dodson
13 R., Cooper L., Van Slyke C.E., Dahdul W., Mabee P., Li D., Harris B., Gillespie
14 M., Jimenez S., Roberts P., Matthews L., Becker K., Drabkin H., Bello S.,
15 Licata L., Chatr-aryamontri A., Schaeffer M.L., Park J., Haendel M., Van
16 Auken K., Li Y., Chan J., Muller H.-M., Cui H., Balhoff J.P., Chi-Yang Wu J., Lu
17 Z., Wei C.-H., Tudor C.O., Raja K., Subramani S., Natarajan J., Cejuela J.M.,
18 Dubey P., Wu C. 2013. An overview of the BioCreative 2012 Workshop
19 Track III: interactive text mining task. Database. 2013:bas056.

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

- 1 Midford P.E., Westerfield M., Vision T.J. 2010. Phenex: Ontological
2 annotation of phenotypic diversity. PLoS One. 5:e10500.
- 3 Balhoff J.P., Dececchi T.A., Mabee P.M., Lapp H. 2014. Presence-absence
4 reasoning for evolutionary phenotypes. phenoday2014.bio-lark.org.
- 5 Barnosky A.D., Matzke N., Tomiya S., Wogan G.O.U., Swartz B., Quental T.B.,
6 Marshall C., McGuire J.L., Lindsey E.L., Maguire K.C., Mersey B., Ferrer E.A.
7 2011. Has the Earth's sixth mass extinction already arrived? Nature.
8 471:51–57.
- 9 Burleigh J.G., Alphonse K., Alverson A.J., Bik H.M., Blank C., Cirranello A.L., Cui
10 H., Daly M., Dietterich T.G., Gasparich G., Irvine J., Julius M., Kaufman S., Law
11 E., Liu J., Moore L., O'Leary M.A., Passarotti M., Ranade S., Simmons N.B.,
12 Stevenson D.W., Thacker R.W., Theriot E.C., Todorovic S., Velazco P.M.,
13 Walls R.L., Wolfe J.M., Yu M. 2013. Next-generation phenomics for the Tree
14 of Life. PLoS currents. 5.
- 15 Cannatella D.C. 1985. A phylogeny of primitive frogs (archaeobatrachians).
16 PhD.-thesis, The University of Kansas.
- 17 Clack J.A., Ahlberg P.E., Blom H., Finney S.M. 2012. A new genus of Devonian
18 tetrapod from North-East Greenland, with new information on the lower jaw
19 of *Ichthyostega*. Palaeontology. 55:73–86.
- 20 Cloutier R., Arratia G. 2004. Early diversification of actinopterygians. In: Arratia
21 G., Wilson M. V. H., Cloutier R., editors. The origin and early radiation of

- 1 vertebrates: Honoring Hans-Peter Schultze. Munich. Verlag Dr. Friedrich Pfeil.
2 p. 217-270.
- 3 Cui H. 2012. CharaParser for fine-grained semantic annotation of organism
4 morphological descriptions. *J. Am. Soc. Inf. Sci. Technol.* 63:738–754.
- 5 Dahdul W.M., Balhoff J.P., Blackburn D.C., Diehl A.D., Haendel M.A., Hall B.K.,
6 Lapp H., Lundberg J.G., Mungall C.J., Ringwald M., Segerdell E., Van Slyke
7 C.E., Vickaryous M.K., Westerfield M., Mabee P.M. 2012. A Unified Anatomy
8 Ontology of the Vertebrate Skeletal System. *PLoS One.* 7:e51070.
- 9 Dahdul W.M., Balhoff J.P., Engeman J., Grande T., Hilton E.J., Kothari C., Lapp H.,
10 Lundberg J.G., Midford P.E., Vision T.J., Westerfield M., Mabee P.M. 2010a.
11 Evolutionary Characters, Phenotypes and Ontologies: Curating Data from the
12 Systematic Biology Literature. *PLoS One.* 5:e10708.
- 13 Dahdul W.M., Lundberg J.G., Midford P.E., Balhoff J.P., Lapp H., Vision T.J.,
14 Haendel M.A., Westerfield M., Mabee P.M. 2010b. The teleost anatomy
15 ontology: anatomical representation for the genomics age. *Syst. Biol.*
16 59:369–383.
- 17 Davis M.C. 2013. The deep homology of the autopod: insights from hox gene
18 regulation. *Integr. Comp. Biol.* 53:224–232.
- 19 Deans A.R., Yoder M.J., Balhoff J.P. 2012. Time to change how we describe
20 biodiversity. *Trends Ecol. Evol.* 27:78–84.
- 21 De Queiroz A., Gatesy J. 2007. The supermatrix approach to systematics.

- 1 Trends Ecol. Evol. 22:34–41.
- 2 Drew B.T., Gazis R., Cabezas P., Swithers K.S., Deng J., Rodriguez R., Katz L.A.,
3 Crandall K.A., Hibbett D.S., Soltis D.E. 2013. Lost branches on the tree of
4 life. PLoS Biol. 11:e1001636.
- 5 Fabre P.H., Rodrigues A., Douzery E.J.P. 2009. Patterns of macroevolution
6 among Primates inferred from a supermatrix of mitochondrial and nuclear
7 DNA. Mol. Phylogenet. Evol. 53:808–825.
- 8 Fabrezi M. 2006. Evolución morfológica en Ceratophryinae (Anura,
9 Neobatrachia). J. Zoolog. Syst. Evol. Res. 44:153–166.
- 10 Gatesy J., Matthee C., DeSalle R., Hayashi C. 2002. Resolution of a
11 supertree/supermatrix paradox. Syst. Biol. 51:652–664.
- 12 Gatesy J., Springer M.S. 2004. A Critique of Matrix Representation with
13 Parsimony Supertrees. In: Binida-Emonds O.R.P. Phylogenetic Supertrees:
14 Combining information to reveal the tree of life. Netherlands. Springer. p.
15 369–388.
- 16 Gkoutos G.V., Green E.C.J., Mallon A.-M., Hancock J.M., Davidson D. 2005. Using
17 ontologies to describe mouse phenotypes. Genome Biol. 6:R8.
- 18 Grafen A. 1988. On the uses of data on lifetime reproductive success. In: Editor:
19 Clutton-Brock T.H., Reproductive success: Studies of individual variation in
20 contrasting breeding systems. Chicago. U. Chicago Press. p. 454-471.

- 1 Haendel M., Balhoff J., Bastian F., Blackburn D., Blake J., Bradford Y., Comte A.,
2 Dahdul W., Dececchi T., Druzinsky R., Hayamizu T., Ibrahim N., Lewis S.,
3 Mabee P., Niknejad A., Robinson-Rechavi M., Sereno P., Mungall C. 2014.
4 Unification of multi-species vertebrate anatomy ontologies for comparative
5 biology in Uberon. *J. Biomed. Semantics*. 5:21.
- 6 Harmon L.J., Baumes J., Hughes C., Soberon J. 2013. *Arbor: Comparative*
7 *Analysis Workflows for the Tree of Life*. *PLoS currents*. 5.
- 8 Harris S.R., Pisani D., Gower D.J., Wilkinson M. 2007. Investigating stagnation in
9 morphological phylogenetics using consensus data. *Syst. Biol.* 56:125–129.
- 10 Hedtke S., Patiny S., Danforth B. 2013. The bee tree of life: a supermatrix
11 approach to apoid phylogeny and biogeography. *BMC Evol. Biol.* 13:138.
- 12 Hejnol A., Obst M., Stamatakis A., Ott M., Rouse G.W., Edgecombe G.D.,
13 Martinez P., Baguñà J., Bailly X., Jondelius U., Wiens M., Müller W.E.G., Seaver
14 E., Wheeler W.C., Martindale M.Q., Giribet G., Dunn C.W. 2009. Assessing the
15 root of bilaterian animals with scalable phylogenomic methods. *Proc. R. Soc.*
16 *B.* 276:4261-4270.
- 17 Hill R.V. 2005. Integration of morphological data sets for phylogenetic analysis
18 of Amniota: the importance of integumentary characters and increased
19 taxonomic sampling. *Syst. Biol.* 54:530–547.
- 20 Hiller M., Schaar B.T., Indjeian V.B., Kingsley D.M., Hagey L.R., Bejerano G. 2012.
21 A ‘forward genomics’ approach links genotype to phenotype using

- 1 independent phenotypic losses among related species. *Cell Rep.* 2:817–823.
- 2 Jeffery J. 2001. Pectoral fins of rhizodontids and the evolution of pectoral
3 appendages in the tetrapod stem-group. *Biol. J. Linn. Soc. Lond.* 74:217-
4 236.
- 5 Johanson Z., Joss J., Boisvert C.A. 2007. Fish fingers: digit homologues in
6 sarcopterygian fish fins. *J. Exp. Zool. (Mol Dev Evol)*. 308B:757-768.
- 7 Kazakov Y., Krötzsch M., Simančík F. 2012. ELK reasoner: Architecture and
8 evaluation. Proceedings of the 1st International Workshop on OWL Reasoner
9 Evaluation (ORE-2012).
- 10 Kazakov Y., Krötzsch M., Simančík F. 2013. The Incredible ELK. *J. Automat.*
11 *Reason.* 53:1–61.
- 12 Kearney M., Clark J.M. 2003. Problems due to missing data in phylogenetic
13 analyses including fossils: a critical review. *J. Vert. Paleontol.* 23:263–274.
- 14 Kearney M. 2002. Fragmentary taxa, missing data, and ambiguity: mistaken
15 assumptions and conclusions. *Syst. Biol.* 51:369–381.
- 16 Kluge, A.G., 1989. A concern for evidence and a phylogenetic hypothesis of
17 relationships among *Epicrates* (Bovidae, Serpentes). *Syst. Zool.* 38:7–25.
- 18 Laurin M., Reisz R. 1997. A new perspective on tetrapod phylogeny. In: Sumida
19 S.S., Martin K.L.M., editors. *Amniote Origins: completing the transition to*
20 *land.* Waltham (MA). Academic Press. p. 9–60.

- 1 Lundberg J.G. 1992. The phylogeny of ictalurid catfishes: A synthesis of recent
2 work. In: Mayden R.L., editor. Systematics, Historical Ecology, & North
3 American Freshwater Fishes. Stanford. Stanford University Press. p. 392–
4 420.
- 5 Mabee P., Balhoff J.P., Dahdul W.M., Lapp H., Midford P.E., Vision T.J.,
6 Westerfield M. 2012. 500,000 fish phenotypes: The new informatics
7 landscape for evolutionary and developmental biology of the vertebrate
8 skeleton. *J. Appl. Ichthyol.* 28:300–305.
- 9 Maddison, W.P. and D.R. Maddison. 2011. Mesquite: a modular system for
10 evolutionary analysis.
- 11 Malia M.J. Jr, Lipscomb D.L., Allard M.W. 2003. The misleading effects of
12 composite taxa in supermatrices. *Mol. Phylogenet. Evol.* 27:522–527.
- 13 Midford P., Dececchi T., Balhoff J., Dahdul W., Ibrahim N., Lapp H., Lundberg J.,
14 Mabee P., Sereno P., Westerfield M., Vision T., Blackburn D. 2013. The
15 vertebrate taxonomy ontology: a framework for reasoning across model
16 organism and species phenotypes. *J. Biomed. Semantics.* 4:34.
- 17 Mungall C.J, Gkoutos G., Washington N., Lewis S. 2007. Representing
18 phenotypes in OWL. OWL: Experiences and Directions (OWLED 2007),
19 Innsbruck, Austria.
- 20 Mungall C.J., Gkoutos G.V., Smith C.L., Haendel M.A., Lewis S.E., Ashburner M.
21 2010. Integrating phenotype ontologies across multiple species. *Genome*
22 *Biol.* 11:R2.

- 1 Mungall C.J., Torniai C., Gkoutos G.V., Lewis S.E., Haendel M.A. 2012. Uberon, an
2 integrative multi-species anatomy ontology. *Genome Biol.* 13:R5.
- 3 Nakagawa S., Freckleton R.P. 2008. Missing inaction: the dangers of ignoring
4 missing data. *Trends Ecol. Evol.* 23:592–596.
- 5 Nakagawa S., Freckleton R.P. 2011. Model averaging, missing data and multiple
6 imputation: a case study for behavioural ecology. *Behav. Ecol. Sociobiol.*
7 65:103–116.
- 8 O’Leary M.A., Bloch J.I., Flynn J.J., Gaudin T.J., Giallombardo A., Giannini N.P.,
9 Goldberg S.L., Kraatz B.P., Luo Z.-X., Meng J., Ni X., Novacek M.J., Perini F.A.,
10 Randall Z.S., Rougier G.W., Sargis E.J., Silcox M.T., Simmons N.B., Spaulding
11 M., Velazco P.M., Weksler M., Wible J.R., Cirranello A.L. 2013. The placental
12 mammal ancestor and the post-K-Pg radiation of placentals. *Science.*
13 339:662–667.
- 14 Price S.A., Wainwright P.C., Bellwood D.R., Kazancioglu E., Collar D.C., Near T.J.
15 2010. Functional innovations and morphological diversification in parrotfish.
16 *Evolution.* 64:3057–3068.
- 17 Ramírez M.J., Coddington J.A., Maddison W.P., Midford P.E., Prendini L., Miller J.,
18 Griswold C.E., Hormiga G., Sierwald P., Scharff N., Benjamin S.P., Wheeler
19 W.C. 2007. Linking of digital images to phylogenetic data matrices using a
20 morphological ontology. *Syst. Biol.* 56:283–294.
- 21 Rudall P.J., Hilton J., Bateman R.M. 2013. Several developmental and

- 1 morphogenetic factors govern the evolution of stomatal patterning in land
2 plants. *New Phytol.* 200:598–614.
- 3 Ruta M., Coates M.I., Quicke D.L.J. 2003. Early tetrapod relationships revisited.
4 *Biol. Rev. Camb. Philos. Soc.* 78:251–345.
- 5 Ruta M. 2011. Phylogenetic signal and character compatibility in the
6 appendicular skeleton of early tetrapods. *Studies on Fossil Tetrapods.*
7 Oxford: Wiley-Blackwell. p. 31–43.
- 8 Sanderson, M.J., 1998. Phylogenetic supertrees: assembling the trees of life.
9 *Trends Ecol. Evol.* 13:105–109.
- 10 Seltmann K.C., Yoder M.J., Mikó I., Forshage M. 2012. A hymenopterists' guide
11 to the Hymenoptera Anatomy Ontology: utility, clarification, and future
12 directions *J. Hymenoptera Res.* 27:67-88.
- 13 Sereno P.C., Tan L., Brusatte S.L., Kriegstein H.J., Zhao X., Cloward K. 2009.
14 Tyrannosaurid skeletal design first evolved at small body size. *Science.*
15 326:418–422.
- 16 Sereno P.C. 2009. Comparative cladistics. *Cladistics.* 25:624–659.
- 17 Shubin N.H., Daeschler E.B., Jenkins F.A. Jr. 2014. Pelvic girdle and fin of
18 *Tiktaalik roseae*. *Proc. Natl. Acad. Sci.* 111:893–899.
- 19 Sigurdson T., Green D.M. 2011. The origin of modern amphibians: a re-
20 evaluation. *Zool. J. Linn. Soc.* 162:457–469.

- 1 Smith B., Ceusters W., Klagges B., Köhler J., Kumar A., Lomax J., Mungall C.,
2 Neuhaus F., Rector A.L., Rosse C. 2005. Relations in biomedical ontologies.
3 Genome Biol. 6:R46.
- 4 Stewart T.A., Smith W.L., Coates M.I. 2014. The origins of adipose fins: an
5 analysis of homoplasy and the serial homology of vertebrate appendages.
6 Proc. R. Soc. B. 281:20133120
- 7 Stoltzfus A., O'Meara B., Whitacre J., Mounce R., Gillespie E.L., Kumar S.,
8 Rosauer D.F., Vos R.A. 2012. Sharing and re-use of phylogenetic trees (and
9 associated data) to facilitate synthesis. BMC Res. Notes. 5:574.
- 10 Swartz B. 2012. A marine stem-tetrapod from the Devonian of western North
11 America. PLoS One. 7:e33683.
- 12 Swenson N.G. 2014. Phylogenetic imputation of plant functional trait databases.
13 Ecography. 37:105–110.
- 14 Van Bocxlaer I., Loader S.P., Roelants K., Biju S.D., Menegon M., Bossuyt F.
15 2010. Gradual adaptation toward a range-expansion phenotype initiated the
16 global radiation of toads. Science. 327:679–682.
- 17 Vos R.A., Balhoff J.P., Caravas J.A., Holder M.T., Lapp H., Maddison W.P., Midford
18 P.E., Priyam A., Sukumaran J., Xia X., Stoltzfus A. 2012. NeXML: rich,
19 extensible, and verifiable representation of comparative data and metadata.
20 Syst. Biol. 61:675–689.
- 21 Wiens J.J., Morrill M.C. 2011. Missing data in phylogenetic analysis: reconciling

1 results from simulations and empirical data. *Syst. Biol.* 60:719–731.

2 Wiens J.J., Tiu J. 2012. Highly incomplete taxa can rescue phylogenetic analyses
3 from the negative impacts of limited taxon sampling. *PLoS One.* 7:e42925.

4 Woltering J.M., Noordermeer D., Leleu M., Duboule D. 2014. Conservation and
5 divergence of regulatory strategies at Hox loci and the origin of tetrapod
6 digits. *PLoS Biol.* 12:e1001773.

7 Zhu M., Yu X., Janvier P. 1999. A primitive fossil fish sheds light on the origin of
8 bony fishes. *Nature.* 397:607–610.

9

10

11 **FIGURE CAPTIONS**

12 Figure 1. Flow chart showing computational steps used to extract synthetic
13 presence/absence supermatrices from ontology-annotated evolutionary
14 phenotypic data. Phenotypic character states of taxa from the evolutionary
15 literature are semantically annotated using anatomy, quality, and taxon
16 ontologies. Using the phenoscape-kb-owl-tools data processing pipeline
17 (<https://github.com/phenoscape/phenoscape-owl-tools>), these phenotypes are
18 reasoned across and deposited into the Phenoscape Knowledgebase. The
19 OntoTrace tool enables a user to generate synthetic presence/absence matrices
20 for specific taxa (here ‘Sarcopterygii’) and particular anatomical entities (here
21 ‘parts of fin or limb’). These matrices, including provenance for each cell, can be

1 viewed in Phenex.

2 Figure 2. Screenshot from Phenex, showing a portion of synthetic supermatrix
3 in Matrix panel (left), synthetic characters in Characters panel (right), and
4 provenance in the new Supporting State Sources panel (below). Here the
5 Supporting State Sources panel display the sources of the character states for
6 the synthetic character 318 'humerus' in *Ichthyostega stensioei*.

7 Figure 3. Ontology-based inference of presence and absence. Direction of
8 arrows indicate the reasoning pathway. Top: The presence of a structure
9 (humerus) is inferred from an assertion to its shape (humerus L-shaped) or a
10 part (entepicondyle of humerus is present). The presence of humerus implies
11 the presence of forelimb skeleton (humerus is part of a forelimb skeleton), a
12 forelimb (forelimb skeleton is part of a forelimb), and thus a forelimb bud
13 (forelimb develops from a forelimb bud). Bottom: In contrast, an assertion to
14 the absence of a humerus does not entail the absence of a forelimb bud,
15 forelimb, or forelimb skeleton; it does entail the absence of its parts
16 (entepicondyle). However, the absence of a forelimb bud entails the absence of
17 a forelimb, thus a forelimb skeleton and thus the humerus.

18 Figure 4. A) Bird's Eye View in Mesquite (Maddison and Maddison 2011)
19 showing inferred (green), asserted (blue), and missing (white) data in the
20 synthetic supermatrix for the first 48 taxa (of 1,051) and all 639 characters. B)
21 Phylogeny of sarcopterygian vertebrates (Tetrapoda in grey) represented in the
22 synthetic supermatrix, showing the distribution of data for one character,

1 'skeleton of digitopodium'. Tip labels in black denote absence of data, blue
2 denotes taxa with asserted data, and green denotes taxa with inferred data.

3 Figure 5. The level of anatomical data available for different parts of the fin and
4 limb can be visualized for taxa along the fin to limb transition. Taxa included in
5 this analysis encompass all major clades from the base of Sarcopterygii to the
6 basal amphibians *Baphetes* and *Westlothiana* (see Ruta 2011 and Clack et al.
7 2012 for source topography). All taxa in this analysis are extinct with exception
8 of the lungfish *Neoceratodus*. Taxa lacking all data for fin or limb were excluded.
9 Cell color reflects the number of character states that entail the presence or
10 absence of that entity for each taxon (row).

11

12

13

14

15

16

17 SUPPLEMENTARY MATERIALS

18 Supplementary Materials Table 1. List of publications used in constructing the
19 synthetic supermatrix. Focal group, number of taxa, and number of fin, limb, and
20 girdle characters, states and phenotype annotations. Studies focused explicitly

1 on the fin to limb transition are denoted by an asterisk.

2 Supplementary Materials Table 2. Taxa (136) present in the variable-only
3 synthetic supermatrix based on inferred data alone.

4 Supplementary Materials Table 3. Conflicting characters. Characters with
5 conflicting states in the variable-only supermatrix, listed by taxon. Conflict type
6 (between direct assertions, direct vs. inferred, and inferred vs. inferred)
7 indicated in right-most column.

8 Supplementary Materials Table 4. Isomorphic characters. Clusters (93) of fully
9 isomorphic characters across the variable-only synthetic supermatrix, arranged
10 from high (10) to low (2)

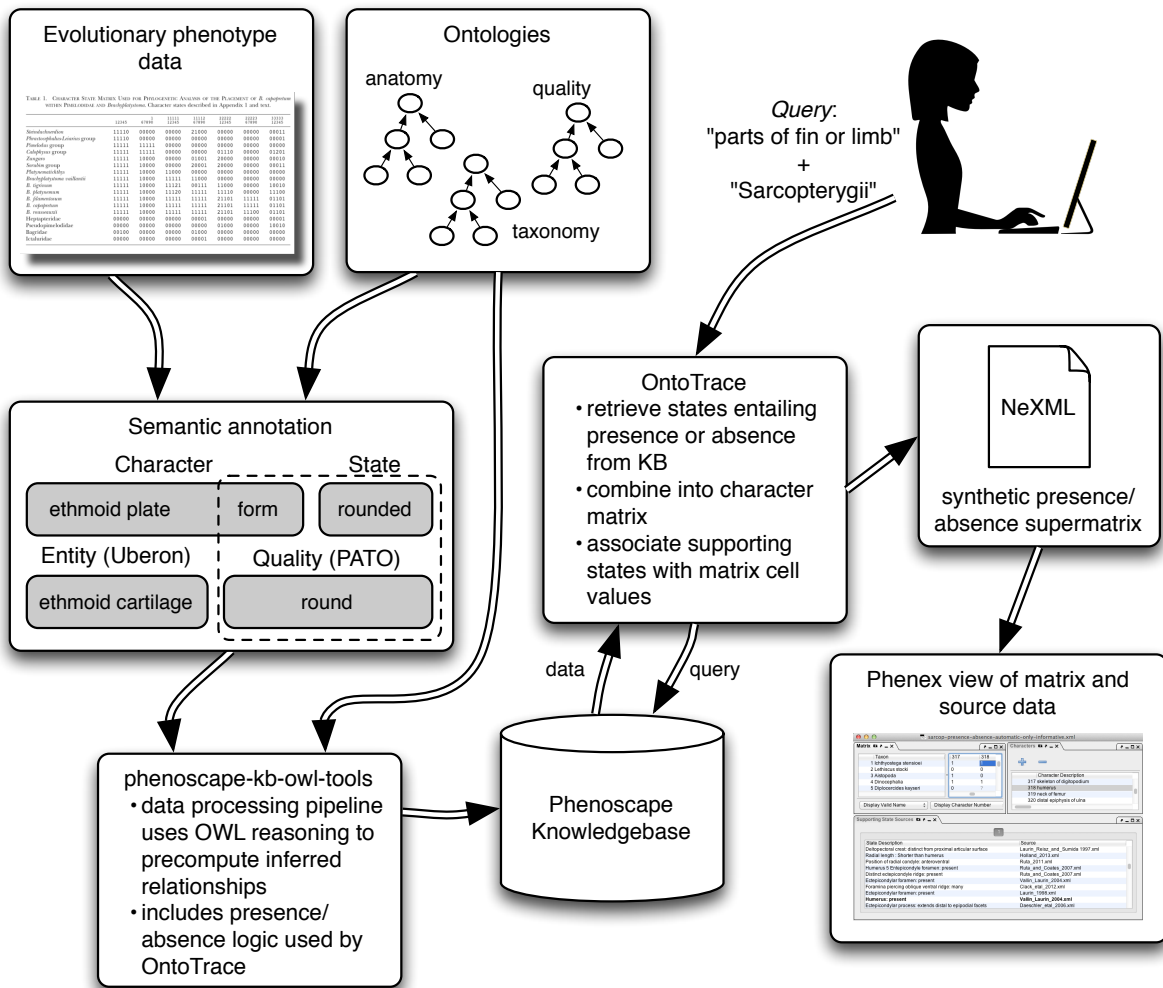
11 Supplementary Materials Table 5. Taxon sampling. The number of source
12 matrices (right-most column) from which taxa (Vertebrate Taxonomy Ontology
13 (VTO) identifier number, left-most column), at various taxonomic ranks, were
14 sampled.

15 Supplementary Materials Table 6. The number of published character states
16 that entail the presence or absence for selected sets of anatomical entities and
17 taxa.

18

19

Figure 1.



PeerJ PrePrints

Figure 2.

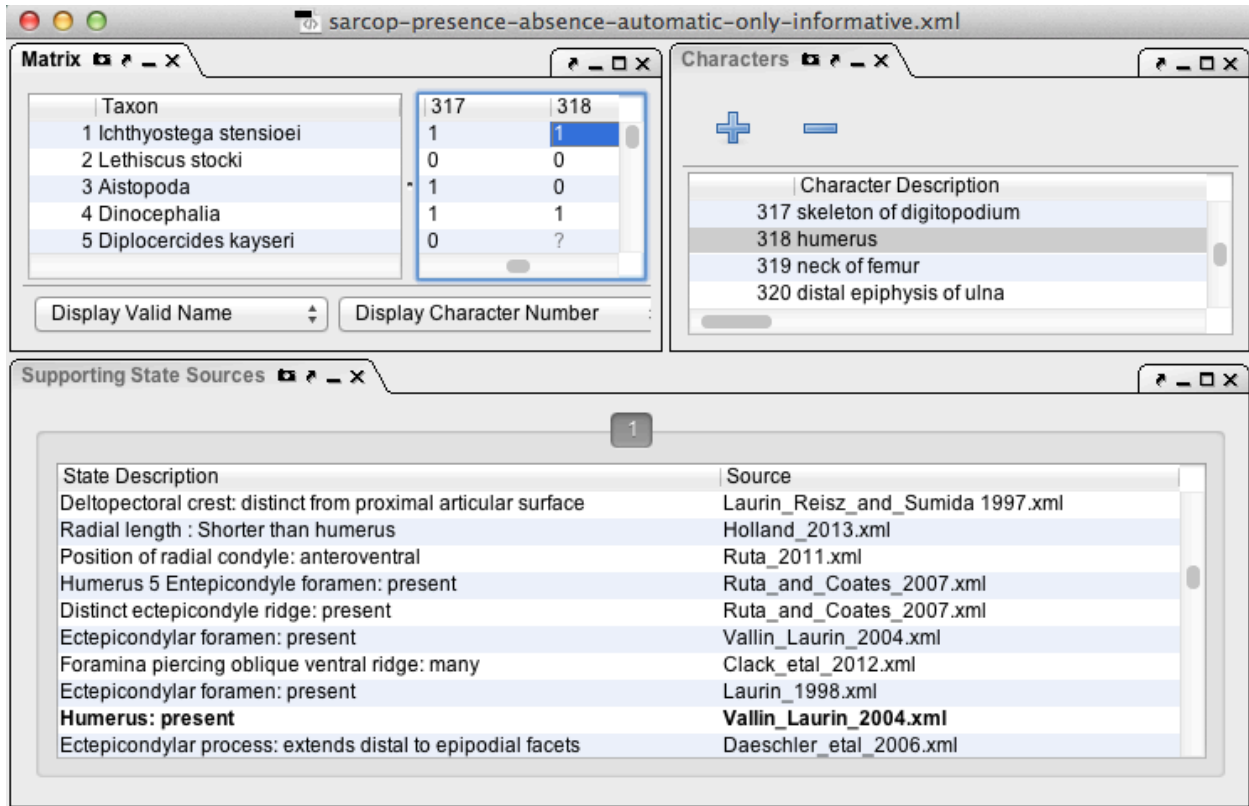


Figure 3.

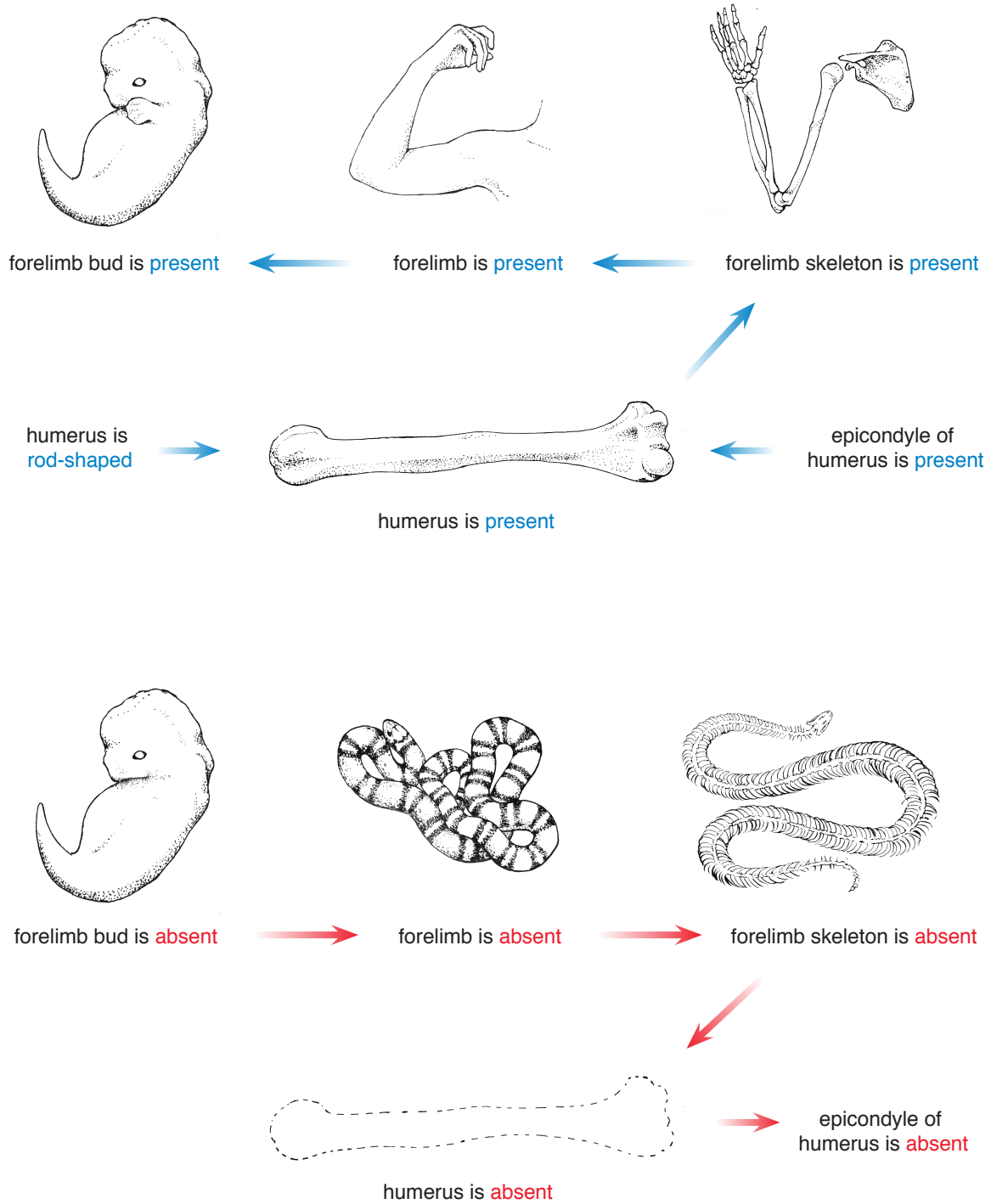


Figure 4A.

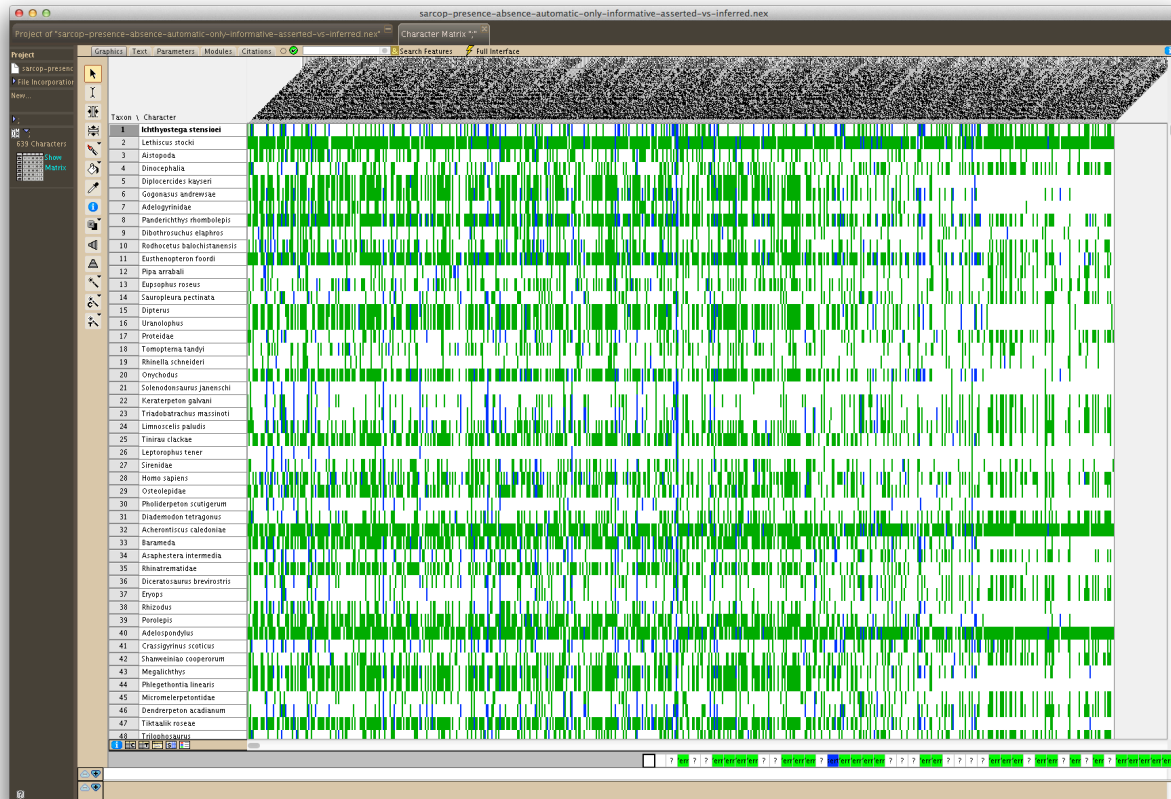


Figure 4B.

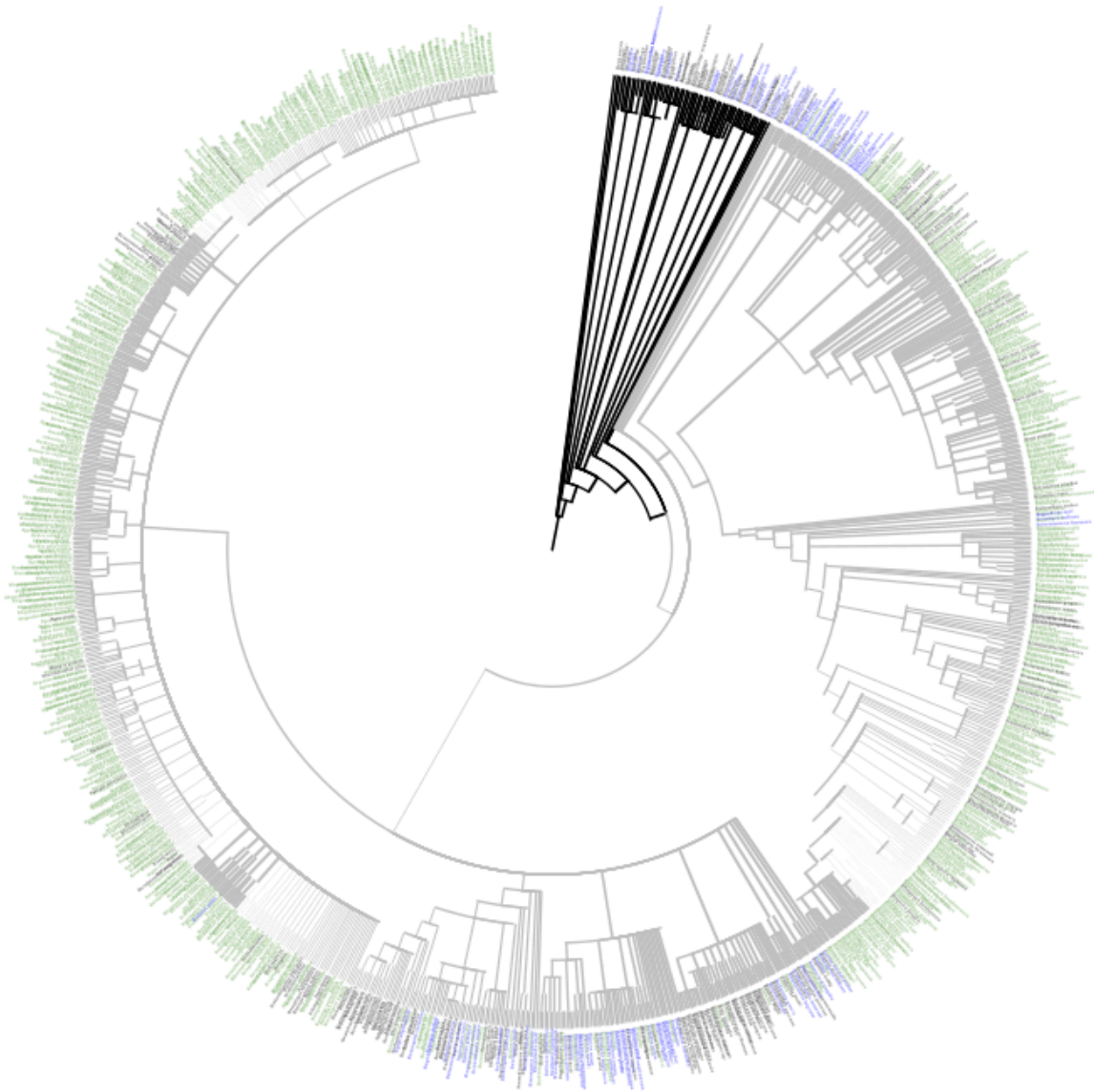


Figure 5.

