

1 *Phyloreferences: Tree-Native, Reproducible, and Machine-Interpretable Taxon Concepts*

2  
3 Nico Cellinese, Florida Museum of Natural History, University of Florida, Gainesville, Florida  
4 32611, U.S.A.

5 Email: [ncellinese@flmnh.ufl.edu](mailto:ncellinese@flmnh.ufl.edu)

6  
7 Stijn Conix, Center for Logic and Philosophy of Science, Katholieke Universiteit, Leuven,  
8 Belgium.

9 Email: [stijn.conix@kuleuven.be](mailto:stijn.conix@kuleuven.be)

10  
11 Hilmar Lapp, Duke Center for Genomics and Computational Biology, Duke University, Durham,  
12 North Carolina 27708, U.S.A.

13 Email: [hilmar.lapp@duke.edu](mailto:hilmar.lapp@duke.edu)

14  
15 **SUBJECT TERMS**

16 Computational semantics

17 Data integration

18 Linnaean names

19 Phylogenetic definitions

20 Phyloreferences

21 Taxon concepts

22 Tree of Life

23  
24 **Abstract**

25 Evolutionary and organismal biology have become inundated with data. At the same rate, we are  
26 experiencing a surge in broader evolutionary and ecological syntheses for which tree-thinking is  
27 the staple for a variety of post-tree analyses. To fully take advantage of this wealth of data to  
28 discover and understand large-scale evolutionary and ecological patterns, computational data  
29 integration, i.e. the use of machines to link data at large scale, is crucial. The most common  
30 shared entity by which evolutionary and ecological data need to be linked is the taxon to which  
31 they belong. We propose a set of requirements that a system for defining such taxa should meet

32 for computational data science: taxon definitions should maintain conceptual consistency, be  
33 reproducible via a known algorithm, be computationally automatable, and be applicable across  
34 the tree of life. We argue that Linnaean names, the most prevalent means of linking data to taxa,  
35 fail to meet these requirements due to fundamental theoretical and practical shortfalls. We argue  
36 that for the purposes of data-integration we should instead use phylogenetic definitions  
37 transformed into formal logic expressions. We call such expressions *phyloreferences*, and argue  
38 that, unlike Linnaean names, they meet all requirements for effective data-integration.

39

## 40 **1. Introduction**

41 The last two decades have witnessed a vast increase of available digital biodiversity data.  
42 This richness in data has been fostered, in part, by a call to mass-digitize museum repositories  
43 (Beaman and Cellinese 2012; Page et al. 2015), and is fueled by the emergence of new  
44 applications and data sources, analytical methods, faster algorithms, and improved  
45 environmental sensors, among others (Philippe et al. 2005; Porter et al. 2009; Michener and  
46 Jones 2012; Chan and Ragan, 2013; Hampton et al. 2017; Kozlov et al. 2019). Additionally, it  
47 has led to a corresponding increasing need for digital access, sharing, and re-purposing of data,  
48 and, consequently, to a need of using machines to link data from different sources to shared  
49 entities. The natural framework for such synthesis of biodiversity data is the Tree of Life. Tree-  
50 thinking has seized a prominent role in systematics since the advent of phylogenetics  
51 (Zimmermann 1931, 1934, 1943; Hennig 1950, 1966). The rapidly increasing knowledge across  
52 the Tree of Life has now enabled a synthesis of phylogenetic hypotheses on a Tree of Life scale,  
53 to produce an encompassing – and digitally fully reusable – view of Life’s evolution, the Open  
54 Tree of Life (Hinchliff et al. 2015; McTavish et al. 2017). As a comprehensive and repeatable  
55 phylogenetic synthesis, it provides unprecedented opportunities for studying evolutionary  
56 patterns across all clades, at large as well as small scales. These clades are the perfect loci at  
57 which to integrate the suite of different data types resulting from evolutionary and biodiversity  
58 research (e.g., Allen et al. 2018; Eliason et al. 2019; Folk et al. 2019; Howard et al. 2019).

59 Thus, a system of defining clades is needed to link the vast amount of available  
60 biodiversity data in a way that it can be recovered, aggregated, and integrated. However, there is  
61 wide disagreement about which system should be used for this purpose. Currently, most  
62 biological data and knowledge are directly or indirectly linked to biological taxa via Linnaean

63 taxon names. As we will discuss below, it is well known that in its current shape the Linnaean  
64 system leads to numerous problems when applied to data-intensive science that depends on  
65 computation. Therefore, an alternative is needed. Broadly speaking, there are two main  
66 candidates for such an alternative: to modify the current Linnaean system such that it can fulfill  
67 certain requirements (see list below), or, more radically, to abandon the Linnaean system in this  
68 context and implement a purely phylogenetic system for clade definitions. The former of these  
69 involves repurposing Linnaean names to refer to clades, and using these names as labels for  
70 taxon concepts<sup>1</sup>. In that sense, this option is a hybrid between the Linnaean and a phylogenetic  
71 system. The latter of these, instead, consists in generating purely phylogenetic definitions of  
72 clades.

73 To arbitrate between these alternatives, we propose the following four requirements that  
74 any system suitable for data-integration should meet: (i) The mapping maintains conceptual  
75 consistency, meaning, when mapped to different phylogenies, the semantics of the retrieved  
76 clades are consistent<sup>2</sup>. (ii) The mapping of a given clade concept to a given phylogenetic  
77 hypothesis is exactly reproducible via a known algorithm. (iii) The algorithm to (re)produce the  
78 mapping is computationally automatable, which is necessary for processing the very large  
79 phylogenies and datasets characteristic of modern biology. This means consulting expert opinion  
80 cannot be part of the algorithm. (iv) The system is applicable to all lineages in the Tree of Life,  
81 including in particular those where Linnaean names are not available (e.g., Archaea, fungi, etc.).

82 In this paper, we show that it is in principle impossible for the Linnaean system to meet  
83 these requirements, and present a purely phylogenetic alternative that does meet them. In section  
84 2 we elaborate on the problems of the Linnaean system, and show that it is beyond repair. In  
85 section 3 we introduce the purely phylogenetic approach, and show how it can address the  
86 shortcomings of the Linnaean system. In section 4 then we introduce one way in which such a  
87 phylogenetic alternative could be implemented, namely, *phyloreferences*, and in section 5 we

---

<sup>1</sup> A taxon concept is the underlying meaning of a group (taxon), whether the group is defined by traits (Linnaean taxonomy) or diagnosed by traits (phylogenetic taxonomy).

<sup>2</sup> By semantics we mean the study, processing, and representation of meaning. The term is used in distinct disciplines, including linguistics and philosophy. In this paper, we use semantics in the sense of computational semantics, which concerns itself with the construction of and automated reasoning with representations of meaning (such as ontologies and logic expressions using ontologies) of natural language expressions.

88 argue that this implementation is preferable over other existing implementations. Finally, in  
89 section 6 we address various objections to our proposal, and section 7 concludes the paper.

90 First of all, it is important to emphasize that the issue at stake in this paper is not that of  
91 nomenclature. The question of how to define taxon concepts for data integration is independent  
92 from the question of whether these taxon concepts also are named, and even whether these  
93 names are Linnaean names. While the approach we propose in this paper fits more naturally with  
94 a form of phylogenetic nomenclature, it is also compatible with retaining Linnaean names.  
95 Related to this, the issue at stake is not that of whether we should recognize certain taxa as  
96 species (Mishler and Wilkins 2018). While a phylogenetic approach like the one proposed here  
97 denies that there is an ontological difference between taxa at different levels, it is compatible  
98 with recognizing some of these taxa as species. Thus, what is at stake is the best way of defining  
99 taxa for data integration, and not the names of these taxa or whether they can be listed as species.

100

101

## 102 **2. The Poverty of Linnaean Names**

103 Many authors before us have pointed to problems caused by Linnaean nomenclature and  
104 classification. This section instead discusses two problems of the Linnaean system that make it  
105 unsuitable for data integration, and argues that it is not possible to eliminate these problems  
106 simply by making small changes to the system.

107

### 108 **2.1. The Linnaean Shortfall Limits Data Discovery**

109 A first problem of the Linnaean system is often referred to as the ‘Linnaean shortfall’.  
110 This is the significant gap in our current knowledge of described vs unknown biological diversity  
111 (Brown and Lomolino 1998; Hortal et al. 2015), and highlights our limited ability to first  
112 discover and then describe taxa according to the rules of nomenclatural codes. In view of the 6<sup>th</sup>  
113 mass extinction we are currently experiencing (Brook et al. 2008), this represents a true plague in  
114 biodiversity science because it implies that we are also losing unknown diversity, and the  
115 diversity we do discover is not described (in a Linnaean framework) fast enough. From a  
116 computational perspective, the latter point represents a true obstacle to addressing the

117 computable taxon concept challenge because taxa need to be described before they can serve as  
118 loci to link data.

119         Two causes of the Linnaean shortfall are particularly relevant in this context. First, the  
120 process of describing diversity is very time consuming and relies on detailed comparative studies  
121 of specimens in museum's repositories and field observations. Second, there are far more levels  
122 of clades in the Tree of Life than there are ranks to name them. As a result, we continue to  
123 discover lineages that persist between revisions of the Tree of Life, yet do not have, and may  
124 never receive, the kind of names required to facilitate discovery and reuse in a name-based  
125 system, let alone formal Linnaean names. Adopted placeholders such as 'phylotype X' or 'clade  
126 A' may serve their purpose within a publication, but they are not discoverable and reusable terms  
127 beyond it (also, see appendix in de Queiroz and Donoghue 2013). This predicament applies  
128 across the Tree of Life, but is particularly prevalent in Archaea and other prokaryotes, and very  
129 common even in many eukaryotes. Consequently, such lineages have often been referred to as  
130 'dark taxa' (Parr et al. 2012).

131         The result is that there are a lot of data about taxa that cannot yet, and may never be,  
132 linked to Linnaean names. This way, the Linnaean system fails to meet requirement (iv), i.e. to  
133 provide the tools to define, communicate and query these unnamed taxa.

134

## 135         2.2. Linnaean names make data discovery difficult to reproduce

136         One might argue that the rate of species descriptions and formal names could, in  
137 principle, increase dramatically and thus alleviate the problem described in the previous  
138 subsection. This subsection argues that even if that were the case, Linnaean names would not be  
139 suitable for integrating data from different sources. This is because it falls short of the three other  
140 requirements as well: (i) it fails to maintain conceptual consistency, (ii) the mapping of a  
141 Linnaean name to a phylogeny is not reproducible by a known algorithm, and (iii) the algorithm  
142 to do this mapping is not automatable.

143         To see why the Linnaean system falls short of these requirements, it is helpful to briefly  
144 consider its design and history. Prior to Linnaeus, biological knowledge was organized in large,  
145 poorly defined categories, and nomenclature was completely unstructured. Linnaeus was a  
146 revolutionary for his time, not so much for the system he created (other botanists before him  
147 experimented with the ranking system), but for what he enabled. He brought order by

148 formalizing criteria to define logical relationships among abstract classes (categorical ranks) and  
149 restructuring the nomenclatural system by enforcing a binomen to every organism at the species  
150 level and a single name to every higher rank. Outside of the – yet to be established – unifying  
151 context of evolution, taxa were assumed to be static entities, with character similarity providing  
152 the best approach to defining groups of organisms. In this context, Linnaean nomenclature  
153 served the need of linking names to taxon groups.

154 Darwinian theory then revolutionized the perspective on biological relationships and  
155 taxon group membership, with the notion that it is natural processes that give rise to taxa, while  
156 characters can only diagnose, but not define categories (Darwin 1859). Zimmermann (1931,  
157 1934, 1943) and Hennig (1950, 1966) formalized these theories and provided the criteria to  
158 construct phylogenetic trees. In this theoretical framework, in which taxa are no longer seen as  
159 static entities, it quickly became clear that the phylogeny-governed hierarchy of Hennig’s  
160 framework is better suited for defining taxa than the logical relatedness of groups in Linnaeus’  
161 hierarchical framework (see also Ereshefsky 2001). Consequently, as common practice Linnaean  
162 nomenclature has been repurposed to link names to clades. In this hybrid system, Linnaean  
163 names are used to label taxon concepts, which are clades rather than fixed entities defined by a  
164 set of characters.

165 However, the Linnaean elements that this hybrid system retains make it impossible to be  
166 used for effective data-integration. There are three reasons for this.

167 First, repurposed Linnaean names define taxon concepts by means of a type specimen  
168 and description (Brzozowski 2020). However, whenever the type is missing from the phylogeny  
169 - which is typically the case - there are no agreed rules for mapping type specimens to clades.  
170 Instead, this mapping relies on expert judgement. As different experts tend to do this in different  
171 ways (see our example of *Campanula* below), this means that the Linnaean system does not meet  
172 requirement (ii) of reproducibility by a single algorithm. In addition, the necessity of expert  
173 judgement means that the mapping of names to clades cannot be automated. This means that the  
174 Linnaean system also fails to meet requirement (iii).

175 Second, the lack of reproducibility in the Linnaean system leads, over time, to confusion  
176 over the taxon concept to which a name is linked. Through time, different experts often apply the

177 same name in different ways due to different interpretations of the original taxon protologue<sup>3</sup>,  
178 and consequently, the meaning of this name becomes difficult to track. This problem is further  
179 exacerbated by purely nomenclatural issues that notoriously plague taxonomy, such as  
180 synonymy, homonymy, misapplication, etc. And even though these can often be reconciled  
181 (albeit not always easily) by taxonomic name resolution services (Boyle et al. 2013; Chamberlain  
182 and Szöcs 2013), this provides little relief to the long-standing informatics challenge of  
183 reconciling names with taxon concepts. This problem is particularly heightened in names with a  
184 long history and legacy of taxonomic literature. Because repurposed Linnaean names still point  
185 to traditionally circumscribed groups that are not generated in an evolutionary framework, they  
186 inherit these problems. In that sense, repurposed Linnaean names approximate to clades, but  
187 never exactly match them. This is because traditional groups and the clades we discover are  
188 fundamentally two different entities, created by very different criteria (Cellinese et al. 2012).  
189 Furthermore, even if the extension of a Linnaean name were to coincide with that of a particular  
190 clade, over time this would quickly fall prey to the same problems of interpretation and  
191 taxonomic as well as phylogenetic revision. Due to the above points, the Linnaean system fails  
192 requirement (i), i.e. it cannot maintain conceptual consistency.

193         Third, the hybrid system still links data to a Linnaean *name*. These names are text strings  
194 without computational meaning. Thus, even if we repurpose a Linnaean name to refer to a clade,  
195 this name can never express the semantics of that clade. Instead of defining the taxon in a way  
196 that would allow machines to identify the taxon, these names link to type specimens and  
197 descriptions that, as described above, have been used and interpreted in different ways by  
198 different researchers. Thus, as long as Linnaean names are used to point to taxon concepts, it will  
199 be impossible for machines to reliably integrate data. This means, again, that the hybrid  
200 Linnaean system inevitably fails to meet the requirement of making taxon definitions  
201 computationally automatable (iii).

202         The failure of the Linnaean system to meet these three requirements is easiest to explain  
203 by drawing an analogy with geolocation-linked data: like taxa, such location data is incredibly  
204 useful for integrating data. Imagine that for geolocation-linked data only place names, not

---

<sup>3</sup> A taxon protologue is the collection of material associated with the publication of a taxon name and concept and therefore, includes all the evidence that support the establishment of a new named entity (e.g., diagnosis, specimens, phylogeny, etc.).

205 standard latitude/longitude geo-coordinates, were available for computation. Data could not be  
206 aggregated by region, users could not draw a bounding box on a map to query a database, species  
207 occurrence data could not be queried for “all species within 50 miles of my location”, and users  
208 querying by place would have to know country, state, and possibly city to make the query less  
209 ambiguous. Yet, this is the current situation in computing with taxon-linked data.

210 Consider, as an example to illustrate the problems of the Linnaean system, the genus  
211 *Campanula* formalized by Linnaeus in 1753, for which *Campanula latifolia* L. was later selected  
212 as a lectotype (Britton and Brown 1913). When discussing *Campanula* L., Lammers (2007)  
213 states that “there is no modern classification which accounts for this large genus in its entirety”  
214 and therefore, the exact number of species is unknown, but the current count is at more than 400.  
215 The original description applied to *Campanula* has been so stretched through time that,  
216 unsurprisingly, *Campanula* as a Linnaean taxon concept is highly polyphyletic, scattered across  
217 the entire Campanuloideae tree with other polyphyletic genera (Crowl et al. 2016; Fig. 1). The  
218 clade including the type specimen (*Campanula latifolia*) would have to retain the original name,  
219 which would imply a cascade of name changes across the tree, not an uncommon repercussion in  
220 taxonomic revisions. Even ignoring the nuisance of name changes, all phylogenetic studies to  
221 date have analyzed a significantly incomplete taxon sample, which had stalled any formal update  
222 in the taxonomy and classification because it would be premature. The most challenging  
223 bottleneck is the inability to retrieve taxonomic concepts unambiguously. Aside from its type  
224 specimen, what constitutes the traditional taxon *Campanula*, in view of how the name has been  
225 applied across time, is not even easy to verbalize, given an author’s subjective taxon description  
226 and the lack of informative synapomorphies. Figure 1 illustrates some of the practical  
227 consequences of this complex issue, by requesting occurrence data from GBIF (gbif.org) using a  
228 query for *Campanula* as a genus. Integrating data obtained in this way with the known  
229 phylogeny will necessarily be very challenging at best, given that *Campanula* as a clade does not  
230 exist.

231 Examples like *Campanula* are very common across all domains at any taxonomic level,  
232 and the harmonization between traditional ideas about life and the phylogenetic approaches we  
233 employ to discover natural entities has become a true impediment to progress in querying,  
234 communicating, and ‘decorating’ all of the parts of the Tree of Life in a consistent and



235 reproducible way. In the next section, we discuss an alternative way of defining taxon concepts  
236 for data integration that does not suffer from the problems of the Linnaean system.

237

### 238 **3. The richness of Phylogenetic Definitions**

239 Starting in the mid 1980's a number of authors suggested that taxon names could be  
240 defined by reference to a part of a phylogenetic tree, prompting an extensive theoretical  
241 discussion, as well as the first attempts to generate phylogenetic definitions (Ghiselin 1984;  
242 Gauthier and Padian 1985; Gauthier 1986; Rowe 1987; de Queiroz 1987, 1988; Gauthier et al.  
243 1988; Estes et al. 1988). A phylogenetic definition represents a formal statement that describes a  
244 clade in a phylogeny. This body of work laid the foundation for phylogenetic taxonomy, later  
245 renamed phylogenetic nomenclature, which takes a strictly tree-thinking approach to biological  
246 nomenclature (de Queiroz and Gauthier 1990, 1992, 1994). Soon thereafter, the *PhyloCode*  
247 ([www.phylocode.org](http://www.phylocode.org)) was drafted as an application of phylogenetic nomenclature's principles.

248 Many systematics papers (e.g., de Queiroz 1992, 1994, 1997; Rowe and Gauthier 1992;  
249 Judd et al. 1993, 1994; Bryant 1996, 1997; Sundberg and Pleijel 1994; Christoffersen 1995;  
250 Schander and Tholleson 1995; Lee 1996, 1998, 2001; Wyss and Meng 1996; Brochu 1997;  
251 Cantino et al. 1997, 2007; Kron 1997; Baum et al. 1998; Eriksson et al. 1998; Härlin and  
252 Sundberg 1998; Hibbett and Donoghue 1998; Alverson et al. 1999; Pleijel 1999; Sereno 1999;  
253 Bremer 2000; Brochu and Sumrall 2001) clearly articulated the need to communicate parts of the  
254 Tree of Life and demonstrated that Life could be described by using three basic clade types and  
255 their associated phylogenetic definitions. These are (1) minimum clade definitions, denoting the  
256 smallest clade that includes the most recent common ancestor, and all its descendants, of two or  
257 more internal specifiers; (2) maximum clade definitions, denoting the largest clade that includes  
258 the first ancestor, and all its descendants, of one or more internal specifiers but excludes one or  
259 more external specifiers; and (3) apomorphy-based definitions, denoting the clade that arises  
260 from the first ancestor, and includes all its descendants, that possesses a specified character that  
261 is synapomorphic with an internal specifier (Fig. 2). Specifiers are reference points in the  
262 phylogeny that serve as anchors for the clade definition and these can be species, specimens, or  
263 apomorphies, which would include molecular sequences. Ideally, when using species as  
264 specifiers, these would already have a phylogenetic definition available or the Linnaean type

265 present in the phylogeny; likewise, when using apomorphies, ideally every trait used as specifier  
266 should be semantically defined.

267         While there has been extensive debate in the literature (Benton 2000; Blackwell 2002;  
268 Schuh 2003; Polaszek and Wilson 2005; Rieppel 2006; Stevens 2006; de Queiroz and Donoghue  
269 2011; among many others) about possible advantages and disadvantages of the PhyloCode as a  
270 nomenclatural system, the PhyloCode is simply one application of phylogenetic nomenclature, in  
271 the realm of nomenclatural codes. Our concern here is not arguing the merits of, or issues with  
272 the PhyloCode, or, for that matter, any nomenclatural code. Instead, we posit that phylogenetic  
273 definitions have unquestionable benefits as a means to unambiguously label all clades in the Tree  
274 of Life, and use these for data integration.

275         Compared to traditional taxon descriptions, phylogenetic definitions have clear  
276 advantages for computing with taxon concepts in a phylogenetic context. They draw  
277 unambiguous reference to any part of the Tree of Life and can be expressed in a formal and  
278 standardized format. Although when published they refer to a taxon concept (clade) originating  
279 from a specific phylogenetic topology, a formal clade concept established by an author is an  
280 unambiguous statement and approach to communicate taxa, and thus data for those taxa,  
281 regardless of future changes in phylogenetic knowledge. That is, as long as the specifiers used in  
282 a clade definition have been matched to a given phylogenetic tree, there is no arguing about the  
283 clade identified by the definition<sup>4</sup>. Obviously, this cannot prevent or resolve disagreements about  
284 the actual taxon concept, but it does enable clearly articulating which element(s) of a  
285 phylogenetic definition is(are) the point(s) of contention. In other words, disagreement over a  
286 concept does not imply ambiguity over what the concept represents. Additionally, a change in  
287 phylogenetic knowledge after the original publication of a phylogenetically defined clade  
288 concept may result in taxa now included in the clade that the original author did not intend to be  
289 included, or for which the community is divided about the merits of their inclusion. Definitions  
290 constructed in some ways will prove more robust, in the judgement of the community, than those  
291 built in other ways. However, whether judged “robust” and agreed upon or not, phylogenetic  
292 definitions will always unambiguously point to the same clade on any tree containing all its  
293 specifiers. For example, our definition of Campanulaceae is “the clade originating with the most

---

<sup>4</sup> We come back to the problem of matching specifiers in section 6.1.

294 recent common ancestor of *Campanula latifolia* Linnaeus and all extant organisms or species  
295 that share a more recent common ancestor with *Campanula latifolia* than with *Roussea simplex*  
296 (*Rousseaceae*) J. E. Smith, *Pentaphragma ellipticum* (*Pentaphragmataceae*) Poulsen, or  
297 *Stylidium graminifolium* (*Stylidiaceae*) Swartz ex Willdenow” (Fig. 3; Cellinese 2020).

298 Others may disagree with this definition, however, there is no ambiguity about the  
299 concept being referred to, and the clade it would identify on a given phylogeny.

300 Phylogenetic definitions are not only beneficial at higher (above species), but also at  
301 shallow (species or below-species) taxonomic levels. For example, reconciling Linnaean names  
302 with polyphyletic taxa, which are very common across all domains of life, is clearly non-trivial.  
303 Often, clades can be diagnosed by interesting morphological or genetic synapomorphies.  
304 Traditional taxon names offer little help in referring to such clades, especially if, as is very  
305 common, type specimens are missing from the analyses. For example, Crowl et al. (2015) found  
306 that *Campanula erinus*, a widespread taxon in the Mediterranean basin, nested in a clade of  
307 narrow Aegean archipelago endemics, is polyphyletic and polyploid. In a more in-depth study,  
308 Crowl et al. (2017) discovered cryptic diversity within this species due to hybridization with *C.*  
309 *creutzburgii*, which revealed a hybrid lineage that is morphologically identical to *C. erinus*, but  
310 differs by having a different ploidy (8x vs the parental 4x). An apomorphy-based clade definition  
311 using the trait octoploidy now allows the semantically unambiguous taxonomic recognition of  
312 this otherwise cryptic group (Crowl and Cellinese 2017).

313 Likewise, in other domains, in particular fungi and bacteria, taxa are often so poorly  
314 known that only unnamed “phylotypes” can be identified (e.g., Massana et al. 2000; Kim et al.  
315 2012; Lin et al. 2014; Hibbett 2016). Phylogenetic definitions can address these cases, because  
316 specifiers can use any uniquely identifiable object suitable for matching the taxonomic unit  
317 represented by nodes in a tree. To illustrate this point, in the above Campanulaceae example, the  
318 taxonomic unit identified by having scientific name *Campanula latifolia* could also be identified  
319 by molecular sequence(s) (e.g., “GenBank: EF141027”), or, as in Crowl and Cellinese (2017),  
320 using a specific herbarium specimen with a globally unique identifier.

321 This potential extends below the species level, for example, to label and query  
322 monophyletic entities corresponding to subsets of populations or polyploid derivatives that show  
323 interesting evolutionary and/or biogeographic patterns, but are currently unnamed. These entities  
324 are not considered ‘species’ and a clear mechanism to name them is lacking from all of the

325 formal nomenclature codes. For data publishing, aggregation, and retrieval systems built around  
326 names instead of meaning, data for such entities cannot be recovered, certainly not  
327 computationally.

328         These advantages of phylogenetic definitions are widely acknowledged, and phylogenetic  
329 definitions have been applied across multiple biological domains in numerous recent  
330 phylogenetic studies, resulting in the publication of many clade names, some of which were  
331 subsequently repurposed in other analyses (Borchiellini et al. 2004; Joyce et al. 2004; Cantino et  
332 al. 2007; Conrad et al. 2011; Soltis et al. 2011; Adl et al. 2012; Cárdenas et al. 2012; Hill et al.  
333 2013; Mannion et al. 2013; Schoch 2013; Sterli et al. 2013; Torres-Carvajal and Mafla-Endara  
334 2013; Wojciechowski 2013; Clemens et al. 2014; Hundt et al. 2014; Rabi et al. 2014; Sferco et  
335 al. 2015; Madzia and Cau 2015; Spatafora et al. 2016; Crowl and Cellinese 2017; Wright et al.  
336 2017; Hibbett et al. 2018; de Queiroz et al. 2020; among numerous others). Arguably, this  
337 constitutes ample evidence that generating and using taxon concepts defined by patterns of  
338 ancestry constitutes an increasing need by the community, and that there is a growing consensus  
339 on how to define and use names for such concepts.

340

341

#### 342         **4. What is a Phyloreference**

343         In the form commonly published by authors, phylogenetic definitions, whether following  
344 strict rules of a nomenclatural code (such as the PhyloCode) or not, are natural language text  
345 expressions. In this form, the ability to compute with the semantics expressed in the text, as  
346 requirement (iii) demands, is severely limited. However, unlike definitions in the Linnaean  
347 system, it is possible to transform phylogenetic definitions in natural language text into  
348 computable representations and thereby make their semantics accessible to machines. We  
349 develop a system for such transformations here, and refer to these computable representations as  
350 *phyloreferences*. Specifically, a phyloreference is a representation of a phylogenetic definition as  
351 a formal, logic expression that makes its semantics explicit and machine-accessible through the  
352 use of terms drawn from ontologies. In this way, phyloreferences are an informatics tool for  
353 communicating taxon concepts to machines, as opposed to, for example, a stand-in for Linnaean  
354 (or other) nomenclature. As an informatics tool, phyloreferences harness the theoretical, as well  
355 as applied, results from a wealth of earlier work in phylogenetic nomenclature to enable

356 machines to integrate and navigate organism-linked data by concepts not afforded by Linnaean  
357 taxonomies.

358         Our proposed approach is based on the Web Ontology Language (OWL 2) (W3C OWL  
359 Working Group 2012) Description Logic (DL) framework. OWL has been widely adopted  
360 across the life sciences for representing domain knowledge in machine-processable form as  
361 ontologies (Mungall et al. 2010, 2011, 2012; Vogt 2009; Jensen and Bork 2010; Deans et al.  
362 2011, 2015; Dahdul et al. 2014; Haendel et al. 2014; Thessen et al. 2015; Senderov et al. 2018).  
363 In the context of information science, in which our approach is based, an ontology is a  
364 representational model of a knowledge domain, specifically the concepts (represented as classes)  
365 comprising the domain, and the relationships that hold between them (represented as  
366 relationships between class members). Ontologies have revolutionized our ability to compute  
367 with the semantics of natural language expressions. For example, by linking terms in free text  
368 phenotype descriptions to formal concepts in community ontologies for the relevant knowledge  
369 domains, machine reasoners and statistical algorithms can be used to compute quantitative  
370 metrics for the semantic similarity of different phenotype descriptions (Pesquita et al. 2009;  
371 Washington et al. 2009; Vision et al. 2011; Bauer et al. 20012; Mabee et al. 2012; Manda et al.  
372 2015; Mabee et al. 2018). Enabling machines to understand the semantics of clade definitions for  
373 the purposes of computational data integration is a much less complex task. Nevertheless, clades  
374 used by researchers to aggregate or communicate data arguably form part of our body of  
375 knowledge about the evolution of the tree of life, and it would thus seem prudent to render it as  
376 much computable as other life science knowledge domains.

377         To afford such capabilities to clade definitions, we propose a model of phyloreferences as  
378 defined OWL classes<sup>5</sup>. In this model, the semantics of a phyloreference, and thus the clade  
379 concept it represents, are declared by a so-called OWL class expression, which essentially gives  
380 the necessary and sufficient conditions for class membership. For a class defined in this way,

---

<sup>5</sup> By class we mean a concept in an ontology, and thus an abstract object (in contrast to individuals or instances, which are concrete objects). Unless stated otherwise, in our use classes have intensional rather than extensional definitions, meaning their descriptions state constraints that must be true for an individual object to be a member of the class. The constraints can be stated in natural language, or as a set of logic conditions. In the latter case, a reasoner can infer class membership. Similarly, we use the term individual in the sense of an individual member of a group. The usage of this term should not be confused with the question of whether taxa are, in a metaphysical sense, classes or individuals. We hold that, depending on the epistemic context, taxa can be construed as both individuals and kinds (see also Brigandt 2009). Hence, the approach we take here is compatible with the view that taxa are, in a metaphysical sense, individuals.

381 software tools called reasoners can (among other things) infer for any individual that all  
382 individuals that fulfill all conditions necessarily must be instances of the class. We then model  
383 the topology of a given phylogeny by declaring its nodes as individuals, and asserting  
384 relationships between those that reflect the topological relationships between nodes. This allows  
385 a reasoner to infer which nodes in the phylogeny, if any, match a given phyloreference. This  
386 class expression-based model also enables other inferences through computational reasoning. For  
387 example, aside from inferring class membership of individuals, OWL reasoners can use these to  
388 infer which phyloreferences are equivalent, and which are subclasses of another. Where found,  
389 such relationships would be implied solely by the semantics of the clade as represented in the  
390 OWL class definition, and as such would hold universally. This is in contrast to approaches that  
391 attempt to map Linnaean names to clades in a tree by comparing the clade on the tree and the  
392 Linnaean taxon concept based on the relationship (inclusion, overlap, etc.) between their  
393 respective sets of members (see “Other Efforts” below).

394 As argued in the large body of work on phylogenetic nomenclature on which we have  
395 based our approach, our proposed models for phyloreference expressions represent patterns of  
396 shared and divergent descent, as included and excluded lineages. To illustrate this, a  
397 phyloreference for the clade Campanuloideae might be expressed in OWL like this (OWL  
398 Manchester Syntax (Horridge and Patel-Schneider 2012); properties in italics; for readability,  
399 ontologies of constituent terms are omitted, and term labels are used in place of identifiers):

400

401 <Campanuloideae> EquivalentTo *includes\_TU* some <Campanula\_latifolia> and *excludes\_TU*  
402 some <Lobelia\_cardinalis>.

403

404 This expression<sup>6</sup> models a maximum clade definition and asserts that the class  
405 Campanuloideae is logically equivalent to the set of nodes that include the taxon concept (TU,  
406 for Taxonomic Unit) ‘Campanula\_latifolia’, and exclude the taxon concept ‘Lobelia\_cardinalis’,  
407 two necessary and sufficient conditions (called property restrictions in OWL). The properties

---

<sup>6</sup>The token “some” in the phyloreference example is from OWL Manchester Syntax and signifies existential quantification. Existential quantification (as opposed to universal quantification) properly represents the semantics of the clade definition: for a taxon concept to be included, *some* instance of it needs to be included, not every possibly existing one (observed or not). Likewise for exclusion. TU here is the class of entities that are instances of a given taxon concept. <Campanula\_latifolia> refers to the TU class, “some <Campanula\_latifolia>” is some instance of that class.

408 *includes\_TU* and *excludes\_TU* are drawn from an ontology, specifically, the Phyloreferencing  
409 Ontology, an application ontology that we are developing on top of the Comparative Data  
410 Analysis Ontology (CDAO) (Prosdocimi et al. 2009) for defining the semantics of clade  
411 definition components. For example, *includes\_TU* as a property is defined such that in the above  
412 definition “*includes\_TU* some <*Campanula\_latifolia*>” is true for all nodes that represent an  
413 instance of the taxon concept *Campanula latifolia*, or from which such a node descends. In  
414 contrast, in the above definition “*excludes\_TU* some <*Lobelia\_cardinalis*>” is true for nodes that  
415 have a sibling node representing an instance of the taxon concept *Lobelia cardinalis*, or from  
416 which such a node descends. The semantics of a definition with these properties are transparent,  
417 unambiguous, and readable by machines. As an ontology class, the definition does not pinpoint  
418 one particular node in one particular taxonomy or phylogeny, but the set of all nodes that satisfy  
419 the definition. Because the definition is a formal logic expression, class membership can be  
420 inferred computationally by a reasoner.

421 By defining phyloreferences as ontology classes, their adoption, reuse, unambiguous  
422 reference, and even community vetting can be promoted using the same mechanisms as for other  
423 widely used community ontologies in the life sciences. Specifically, they can be given a label,  
424 allowing reference to them by name; assigned globally unique identifiers, making them  
425 unambiguously referenceable; and assembled into an ontology maintained in an infrastructure,  
426 such as a Github repository that facilitates version control, releases, and community  
427 collaboration.

428 Ultimately, a phyloreference in our approach bears the following important properties.  
429 Foremost, it meets our four requirements. Its semantics are unambiguous and machine  
430 interpretable because they are expressed in formal logics with uniquely identified ontology  
431 terms. This enables reproducing their mapping to a given phylogeny with a fully computational  
432 algorithm (requirements (ii) and (iii), and enables maintaining semantic consistency when  
433 mapped to different (such as updated) phylogenies (requirement (i)). When a phyloreference is  
434 applied to a particular phylogeny that lacks a clade with consistent semantics, there will not be a  
435 node that “matches” (i.e. can be inferred as an instance). As a logically defined ontology class, a  
436 phyloreference can but need not be named. If it is named, the name is only a label to aid human  
437 communication, and this label does not carry semantics a machine is expected to recognize.  
438 Phyloreferencing can thus be applied to any branch of the Tree of Life, whether useful names

439 exist or not (requirement (iv)). A phylolreference class can be given a globally unique identifier  
440 by which to unambiguously reference it for machines, independent of whether it has a label.

441 Furthermore, in this way phylolreferences are quite similar to terms in other community  
442 ontologies, and our system therefore interoperates naturally with the communities of practice and  
443 tool ecosystems that have developed around collections of ontologies in different domains, in  
444 particular in the life sciences (Smith et al. 2007).

445

## 446 **5. Other Efforts to improve the computability of taxon concepts**

447 Even though there has been much controversy over the application of phylogenetic  
448 nomenclature (Benton 2000; Blackwell 2002; Schuh 2003; Polaszek and Wilson 2005; Rieppel  
449 2006; Stevens 2006; de Queiroz and Donoghue 2013; among many others), its potential to define  
450 taxon concept semantics in a logical manner with unambiguously expressible meaning has been  
451 recognized before. Hibbet et al. (2005), Keesey (2007), and in part Sereno (2005) and Sereno et  
452 al. (2005), already envisioned mechanisms and applications that would leverage computable  
453 clade definitions to unambiguously retrieve taxa based on shared descent-based specifications.  
454 Keesey (2007) includes a notation and formalism for defining clade names based on  
455 mathematical set theory and operators, using the Mathematical Markup Language (MathML), an  
456 XML derivative, and extensions to it. Keesey's approach, unlike ours, also supports group  
457 concepts that are not monophyletic. However, because MathML is a structured syntax language,  
458 not a formal logic, Keesey's approach requires defining custom, bespoke semantics for his  
459 notations. It also does not lend itself to publishing clade definitions in the form of ontologies that  
460 are readily interoperable with the wealth of other community ontologies increasingly widely used  
461 in biology, and the software support even for only reading and interpreting MathML is limited.  
462 In practice, Keesey's proposal has not been adopted.

463 Thau and Ludäscher (2007) and Thau et al. (2008) proposed to use Region Connection  
464 Calculus (RCC, specifically RCC-5; Randell et al 1992) as a formal logic for computationally  
465 reconciling different Linnaean taxonomies (or taxonomic checklists derived from such  
466 taxonomies) with each other. RCC-5 defines five basic relationships between two entities:  
467 equality, proper inclusion, inverse proper inclusion, overlap, and disjointness. In their approach,  
468 human experts assert which relationship(s), called articulations, hold between the concepts from  
469 different input taxonomies, such as concepts with identical names, or names that exist in only



470 some of the input taxonomies. Experts also assign or relax a number of so-called global (or  
471 latent) taxonomic constraints, such as disjointness of sibling taxa, and parent taxon coverage  
472 (every member of a parent taxon is a member of some child taxon). Thau et al. (2008) show that  
473 certain machine reasoners can prove the consistency (or inconsistency) of different taxonomies  
474 under the asserted articulations and constraints, and can infer minimally informative  
475 relationships (a disjunction of one or more of the RCC-5 base relationships) between concepts.

476 More recently, Franz et al. (2016, 2019) and Cheng et al. (2017) applied this approach to  
477 a variety of complex biological use cases, and also extended it to the challenge of reconciling  
478 concepts from traditional Linnaean nomenclature with clades in a phylogenetic tree, as well as  
479 aligning clade concepts from competing phylogenetic hypotheses. Although evidently useful for  
480 the problem of computationally reconciling taxon concepts, for each new input taxonomy or  
481 phylogenetic hypothesis to be reconciled, a considerable amount of effort from trained human  
482 experts is necessary to create the articulations and constraints, and the resulting assertions still do  
483 not disambiguate or make computable the original intensional semantics of a taxon concept.  
484 Therefore, it does not make the exercise of repurposing Linnaean names for clades in a  
485 phylogenetic tree a less subjective and manual approximation than it necessarily is, because the  
486 concepts at hand are fundamentally different in nature.

## 487 **6. Challenges and Limitations**

488 Previous proposals to replace the Linnaean system with a purely phylogenetic alternative  
489 have proven to be very controversial. As our proposal does not concern taxonomic nomenclature  
490 or classification, many of these controversies are not directly relevant. However, there are  
491 various ways in which opponents might object against the arguments in this paper. We respond  
492 to these briefly, and point to limitations and challenges for our approach.

493

### 494 **6.1. Specifiers**

495 One of the greater challenges in applying phyloreferences on a larger scale, and across  
496 different phylogenetic trees, is that phylogenetic definitions are “anchored” by the specifiers  
497 designating the taxon concepts that are to be included or excluded. Therefore, resolving a  
498 phyloreference on a tree necessarily requires that the anchoring taxon concepts of a  
499 phyloreference, and the taxon concepts linked to (typically terminal) nodes in a phylogeny, can  
500 be “matched” by a reasoner. More specifically, these taxon concepts need to be defined such

501 that the reasoner can infer when a taxon concept used in the phyloreference is congruent with, or  
502 includes, a taxon concept linked to a tree node. In some cases such a match will be exact and  
503 unambiguous, for example, if the specifier and node-linked taxon concept are referenced to the  
504 same globally unique identifier. In practice, matching specifiers between phyloreference and  
505 phylogeny is an inherently non-trivial problem, and matches will range from unambiguous to  
506 approximate. For example, if taxon concept references are, as will commonly be the case,  
507 Linnaean taxon names, even an exact match is not necessarily free of ambiguity, such as when  
508 the names are not demonstrably drawn from the same taxonomy. Indeed, this is the taxonomic  
509 name resolution problem that arises whenever Linnaean taxon names must be reconciled, and the  
510 confidence in name matches will follow the familiar spectrum. Especially for phylogenies with  
511 incomplete taxon sampling, a taxon concept used as specifier in a phyloreference may also be  
512 altogether absent from a tree. The question is, then, whether or not one of the taxon concepts  
513 present on the tree can substitute for the specifier without changing the semantics of the clade  
514 definition. Whether this is possible or not will in turn depend on the definition of the clade and  
515 the phylogeny at hand on which it is to be recovered, and may require sophisticated algorithms to  
516 determine.

517         Phyloreferences by themselves do not obviate the need to match or reconcile Linnaean  
518 taxon names. However, this is due to the prevailing practice of identifying taxon concepts  
519 through names, rather than a specific weakness in the phyloreferencing approach; and because  
520 phyloreferences are in essence uniquely identifiable ontology terms, this problem and the  
521 ambiguity it confers are not re-introduced every time data are linked to a taxon. Furthermore,  
522 how and why a taxon concept for a specifier matches one for a node in a tree can be expressed  
523 through formal axioms in the same logic framework (i.e., OWL2 in our case), and thus be  
524 documented in a fully reproducible manner. For example, if a target phylogeny lacks a node for  
525 *Campanula latifolia*, but contains a node for *Campanula*, a “mapping” axiom asserting that the  
526 concept *Campanula* includes *Campanula latifolia* will allow matching a phyloreference for the  
527 Campanuloideae clade that references *Campanula latifolia* as a specifier that must be included.

528         Finally, it is worth emphasizing that the ambiguity inherent in reconciling names by itself  
529 does not introduce ambiguity into the semantics of the clade definition, though it does render  
530 *recovering* the clade semantics on phylogenies, other than the one used by the original author,  
531 prone to the same problems that beset taxon name matching in general. Creating mapping

532 axioms in an effective and scalable manner may be non-trivial, but we are confident that  
533 solutions to address this challenge can and will be developed. In the meantime, the Open Tree of  
534 Life offers a comprehensive, even if synthetic, phylogeny that is continuously updated with  
535 evolving phylogenetic knowledge, and with names for terminal nodes sourced from dozens of  
536 taxonomies (Rees and Cranston 2017).

537

## 538 6.2. Genealogical discordance

539 It is well-known that, due to phenomena such as lateral gene transfer, hybridization,  
540 introgression, and others, evolution is often not tree-like across all domains of life, including  
541 Archaea, bacteria and fungi. One might worry then that the phyloreferences proposed here are  
542 not suitable for capturing groups whose evolutionary relations are more suitably represented by a  
543 network than by a bifurcating pattern. Although phylogenies are hierarchical, with clades that are  
544 either nested or mutually exclusive, reticulation due to different biological processes results in  
545 partially overlapping clades, with hybrid lineages belonging to both parental clades. Partially  
546 overlapping clades can, in fact, be phylogenetically defined, which demonstrates the flexibility  
547 of this approach. For example, Crowl and Cellinese (2017) illustrate how phylogenetic  
548 definitions apply to lineages derived from hybridization and polyploidy (using ploidy in an  
549 apomorphy-based definition), and allow the naming of cryptic diversity.

550 Phylogenetic reconstructions may generate discordant hypotheses that are best  
551 synthesized by networks rather than bifurcating patterns. For considering the question whether  
552 phyloreferences can be meaningfully applied to such networks, note that in principle the key  
553 concepts used in our approach for encoding the semantics of a clade definition, namely ancestors  
554 and descendants, and taxon concepts included in or excluded from a line of descendents, still  
555 fully apply in networks. Hence, there is no theoretical or technical reason that would prevent  
556 resolving a phyloreference on a phylogenetic network. Nonetheless, a clade retrieved in this way  
557 should be treated with great caution, because at least for now the underlying clade definition will  
558 have almost universally been erected based on a phylogenetic tree, not a network. Therefore, the  
559 benefit of applying phyloreferences to networks as part of, for example, a data integration  
560 project, seems questionable at best.

561

## 562 6.3. Adoption cost

563 One could object that even if phylotaxonomies are in principle preferable over Linnaean  
564 names for integrating data, the cost of adoption would be very high, or high enough to outweigh  
565 the benefits. For a response, we note but set aside the fact that such an argument would attribute  
566 limited value to the problems caused by using the Linnaean system; we disagree that  
567 irreproducible science has only limited costs. Nonetheless, we acknowledge that as for any novel  
568 system for indexing data, for a resource such as GBIF, with huge amounts of data that need to be  
569 queryable very efficiently by a large user community, to fully support phylotaxonomy would  
570 likely have a significant engineering cost. This notwithstanding, we find it important to note that  
571 phylotaxonomies can already be taken advantage of right now, including for data integration  
572 projects, by tapping into and combining already existing technologies. To sketch out an example,  
573 the programming interface (API) to the Open Tree of Life includes a most recent common  
574 ancestor query service that depending on the input parameters returns the common ancestor node  
575 semantically fully consistent with minimum clade and maximum clade definitions, respectively,  
576 that underlie phylotaxonomies. Additional Open Tree of Life query services can then be used to  
577 obtain the species contained by the clade resolved in the previous step, which then in turn allow  
578 querying a database indexed by Linnaean names for data associated with the clade. This  
579 approach can already be used, for example, to find how phylogenetic vs Linnaean names can  
580 result in different inferences, such as geographical distribution.

581

## 582 **7. Final remarks**

583 We strongly believe we are at a crossroad where the idiosyncratic applications of  
584 Linnaean nomenclature and taxonomy to the approach we use to discover and name taxa is  
585 simply untenable in the age of computationally-driven science. Linnaean names represent an  
586 incurable theoretical and practical shortfall (see Sterner and Franz 2017). We suggest that  
587 phylotaxonomy lays the foundation for an informatics infrastructure that enables using the Tree  
588 of Life to organize, query, and navigate our knowledge of biodiversity. Building this foundation  
589 now is timely. Large phylogenies encompassing diverse groups across the tree of life are  
590 published in increasing numbers (e.g., Smith et al. 2011; Hinchliff et al. 2015; Smith and Brown  
591 2018; Howard et al. 2019). Especially for large tree synthesis projects, the need for  
592 phylotaxonomy has already arisen, because it is the basis for persistently and reproducibly  
593 linking data and metadata to internal nodes (i.e. clades) in the tree. There are also parts of the

594 Tree of Life for which a stunning organismal and trait diversity is only just beginning to be  
595 characterized, and for which the traditional fallback of Linnaean names is hardly available, and  
596 perhaps never will be (e.g., microbial diversity, and population-level diversity). Yet, the ability  
597 to unambiguously refer to these groups is necessary, not least to organize, query, and retrieve our  
598 knowledge about any group of interest. In contrast to Linnaean names, phylogenetic definitions  
599 can be created using any identifiable object, including specimens, samples, and sequences. If  
600 appropriately labeled and distributed in community-vetted ontologies, phyloreferences can  
601 provide names and concepts that allow researchers to communicate data and knowledge about  
602 their groups, yet also have fully computable and thus reproducible semantics built-in.

603         One of the key goals of phyloreferences is to enable computationally querying,  
604 navigating, integrating, and visualizing any data linked to groups of organisms, in a way that is  
605 driven by evolutionary relatedness. We have argued that merely repurposing Linnaean names  
606 onto trees cannot achieve this goal. Phyloreferences allow us to compare parts of the Tree of Life  
607 about which we would otherwise not be able to communicate. Consequently, the number of  
608 phylogenetic taxon definitions being published has already increased rapidly in recent years  
609 across multiple domains, signifying that phylogenetic approaches to diagnose taxonomic groups  
610 and their names are being increasingly widely adopted and ideally, every clade discovered  
611 should bear a definition. When translated into formal phyloreferences, the semantics of these  
612 definitions not only become fully accessible to machines, but by curating them into a community  
613 ontology, they become much more findable and reusable compared to when buried in the text of  
614 publications.

615         We believe that a phylogenetic data synthesis encompasses far more than a challenging  
616 topological synthesis. The approach we propose is native to tree-thinking and completely flexible  
617 because phyloreferences adapt seamlessly to changes in phylogenetic knowledge and would  
618 therefore apply to small and large topologies and syntheses. In view of the upcoming publication  
619 of the PhyloCode and the ever-increasing number of published phylogenetic definitions, now is  
620 the time to envision the Tree of Life as a navigable map where clade definitions (taxon concepts)  
621 serve as physical addresses and phyloreferences provide the means to achieve a retraceable  
622 navigation.

623

624 **Acknowledgements**

625 This work was supported by a grant from the National Science Foundation (DBI-1458604 and  
626 DBI-1458484 to Nico Cellinese and Hilmar Lapp, respectively) for which we are very grateful.  
627 Our thanks go to Michael Donoghue, David Hibbett, Pam Soltis and Doug Soltis for their  
628 insightful feedback on an early version of this paper. We also thank Andy Crawl and Gaurav  
629 Vaidya for assistance with figure 2. Additionally, we are very grateful to Dr. Ken and Linda  
630 McGurn who provided a generous gift to assist with the phyloreferencing project's needs.  
631 Finally, we thank Christopher Eliot, Joyce Havstad, and three anonymous reviewers for  
632 assistance and constructive comments that greatly improved our paper.

633

### 634 **References**

635 Adl, S.M., A.G.B. Simpson, C.E. Lane, J. Lukeš, D. Bass, S.S. Bowser, M.W. Brown, et al.

636 2012. "The Revised Classification of Eukaryotes." *Journal of Eukaryotic Microbiology* 59:  
637 429–493.

638

639 Allen, J.M., R.A. Folk, P.S. Soltis, D.E. Soltis, and R.P. Guralnick 2018. "Biodiversity synthesis  
640 across the green branches of the tree of life." *Nature Plants* 5: 11-13.

641 Alverson, W.S., B.A. Whitlock, R. Nyffeler, C. Bayer, and D.A. Baum 1999. "Phylogeny of the  
642 core Malvales: evidence from *ndhF* sequence data." *American Journal of Botany* 86: 1474–  
643 1486.

644

645 Bauer, S., S. Köhler, M.H. Schulz, and P.N. Robinson. 2012. "Bayesian ontology querying for  
646 accurate and noise-tolerant semantic searches." *Bioinformatics* 28: 2502–2508.

647

648 Baum, D.A., W.S. Alverson, and R. Nyffeler. 1998. "A durian by any other name: taxonomy and  
649 nomenclature of the core Malvales." *Harvard Papers in Botany* 3: 315–330.

650

651 Beaman, R.S., and N. Cellinese. 2012. "Mass digitization of scientific collections: New  
652 opportunities to transform the use of biological specimens and underwrite biodiversity  
653 science." *ZooKeys* 209: 7–17.

654

655 Benton, M.J. 2000. “Stems, nodes, crown clades, and rank-free lists: is Linnaeus dead?”  
656 *Biological Reviews* 75: 633–648.  
657

658 Blackwell, J.H. 2002. “One-hundred-year code déjà vu?” *Taxon* 51: 151–154.  
659

660 Borchiellini, C., C. Chombard, M. Manuel, E. Alivon, J. Vacelet, and N. Boury-Esnault. 2004.  
661 “Molecular phylogeny of Demospongiae: implications for classification and scenarios of  
662 character evolution.” *Molecular Phylogenetics and Evolution* 32: 823–837.  
663

664 Boyle, B., N. Hopkins, Z. Lu, J.A.R. Garay, D. Mozzherin, T. Rees, N. Matasci, et al. 2013.  
665 “The taxonomic name resolution service: an online tool for automated standardization of  
666 plant names.” *BMC Bioinformatics*. 14: 16.  
667

668 Bremer, K. 2000. “Phylogenetic nomenclature and the new ordinal system of the angiosperms.”  
669 In *Plant systematics for the 21st century*, edited by B. Nordenstam, G. El Ghazaly, and M.  
670 Kassas, 125–133. London, United Kingdom: Portland Press.  
671

672 Brigandt, I. 2009. “Natural kinds in evolution and systematics: Metaphysical and  
673 epistemological considerations.” *Acta Biotheoretica* 57: 77–97.  
674

675 Britton, C.E., and A. Brown 1913. *An illustrated flora of the Northern United States, Canada*  
676 *and the British Possessions*. New York: Charles Scribner's Sons.  
677

678 Brochu, C.A. 1997. “Synonymy, redundancy, and the name of the crocodile stem-group.”  
679 *Journal of Vertebrate Paleontology* 17: 448–449.  
680

681 Brochu, C.A., and C.D. Sumrall. 2001. “Phylogenetic nomenclature and paleontology.” *Journal*  
682 *of Vertebrate Paleontology* 75: 754–757.  
683

684 Brook, B.W., N.S. Sodhi, and C.J.A. Bradshaw. 2008. “Synergies among extinction drivers  
685 under global change.” *Trends in Ecology and Evolution* 23: 453–460.

686  
687 Brown, J.H., and M.V. Lomolino. 1998. *Biogeography*. Sunderland, MA: Sinauer.  
688  
689 Bryant, H.N. 1996. “Explicitness, Stability, and Universality in the Phylogenetic Definition and  
690 Usage of Taxon Names: A Case Study of the Phylogenetic Taxonomy of the Carnivora  
691 (Mammalia).” *Systematic Biology* 45: 174–189.  
692  
693 Bryant, H.N. 1997. “Cladistic information in phylogenetic definitions and designated  
694 phylogenetic contexts for the use of taxon names.” *Biological Journal of the Linnaean  
695 Society* 62: 495–503.  
696  
697 Brzozowski, J.A. 2020. “Biological taxon names are descriptive names.” *History and  
698 Philosophy of the Life Sciences* 42: 29 <https://doi.org/10.1007/s40656-020-00322-1>  
699  
700 Cantino, P.D., R.G. Olmstead, and S.J. Wagstaff. 1997. “A Comparison of Phylogenetic  
701 Nomenclature with the Current System: A Botanical Case Study.” *Systematic Biology* 46:  
702 313–331.  
703 Cantino, P.D., J.A. Doyle, S.W. Graham, W.S. Judd, R.G. Olmstead, D.E. Soltis, P.S. Soltis, and  
704 M.J. Donoghue. 2007. “Towards a phylogenetic nomenclature of Tracheophyta.” *Taxon* 56:  
705 822-846.  
706  
707 Cárdenas, P., T. Pérez, and N. Boury-Esnault. 2012. “Sponge systematics facing new  
708 challenges.” In *Advances in Sponge Science: Phylogeny, Systematics, Ecology*, edited by  
709 M.A. Becerro, M.J. Uriz, M. Maldonado, and X. Turon, 79–209. London, United Kingdom:  
710 Academic Press.  
711  
712 Cellinese, N., D.A. Baum, and B.D. Mishler. 2012. “Species and Phylogenetic Nomenclature.”  
713 *Systematic Biology* 61: 885–891.  
714  
715 Cellinese, N. 2020. “Campanulaceae.” In *Phylonyms: a Companion to the PhyloCode*, edited by  
716 K. de Queiroz, P.D. Cantino, and J. Gauthier, 381-383. Boca Raton, FL: CRC Press.



717  
718 Chamberlain, S.A., and E. Szöcs. 2013. “Taxize: taxonomic search and retrieval in R.”  
719 *F10000Res* 2: 191.  
720  
721 Chan, C.X., and M.A. Ragan. 2013. “Next-generation phylogenomics.” *Biology Direct* 8: 3.  
722  
723 Cheng, Y.Y., N. Franz, J. Schneider, S. Yu, T. Rodenhäuser, and B. Ludäscher. 2017. “Agreeing  
724 to disagree: Reconciling conflicting taxonomic views using a logic-based approach.”  
725 *Proceeding of the Association for Information Science and Technology* 54: 46–56.  
726  
727 Christoffersen, M.L. 1995. “Cladistic Taxonomy, Phylogenetic Systematics, and Evolutionary  
728 Ranking.” *Systematic Biology* 44: 440–454.  
729  
730 Clemens, W.L., M. Arakaki, P.W. Sweeney, E.J. Edwards, and M.J. Donoghue. 2014. “A  
731 chloroplast tree for *Viburnum* (Adoxaceae) and its implications for phylogenetic  
732 classification and character evolution.” *American Journal of Botany* 101: 1029–1049.  
733  
734 Conrad, J.L., J.C. Ast, S. Montanari, and M.A. Norell. 2011. “A combined evidence  
735 phylogenetic analysis of Anguimorpha (Reptilia: Squamata).” *Cladistics* 27: 230–277.  
736  
737 Crowl, A.A., C. Visger, G. Mansion, R. Hand, H.-H. Wu, G. Kamari, D. Phitos, and N.  
738 Cellinese. 2015. “Evolution and Biogeography of the Endemic *Roucelia* complex  
739 (Campanulaceae: Campanula) in the Eastern Mediterranean.” *Ecology and Evolution*  
740 10.1002/ece3.179.  
741  
742 Crowl, A.A., N. Miles, C. Visger, K. Hansen, T. Ayers, R. Haberle, and N. Cellinese. 2016. “A  
743 global perspective on Campanulaceae: biogeographic, genomic, and floral evolution.”  
744 *American Journal of Botany* 103: 233–245.  
745

746 Crowl, A.A., C. Myers, and N. Cellinese. 2017. "Embracing discordance: Phylogenomic  
747 analyses provide evidence for allopolyploidy leading to cryptic diversity in a Mediterranean  
748 *Campanula* (Campanulaceae) clade." *Evolution* 71: 913–922.

749 Crowl A.A., and N. Cellinese. 2017. "Naming diversity in an evolutionary context: Phylogenetic  
750 definitions of the *Roucela* clade (Campanulaceae/Campanuloideae) and the cryptic taxa  
751 within." *Ecology and Evolution* 10.1002/ece3.3442.

752

753 Dahdul, W.M., H. Cui, P.M. Mabee, C.J. Mungall, D. Osumi-Sutherland, R.L. Walls, and M.A.  
754 Haendel. 2014. "Nose to tail, roots to shoots: spatial descriptors for phenotypic diversity in  
755 the Biological Spatial Ontology." *Journal of Biomedical Semantics* 5: 34.

756

757 Darwin, C. 1859. *On the Origin of Species*. London, United Kingdom: John Murray.

758

759 de Queiroz, K. 1987. *Phylogenetic systematics of iguanine lizards. A comparative osteological*  
760 *study*. Berkeley, CA: University of California Press.

761

762 de Queiroz, K. 1988. "Systematics and the Darwinian Revolution." *Philosophy of Science* 55:  
763 238–259.

764

765 de Queiroz, K. 1992. "Phylogenetic definitions and taxonomic philosophy." *Biology and*  
766 *Philosophy* 7: 295–313.

767

768 de Queiroz, K. 1994. "Replacement of an Essentialistic Perspective on Taxonomic Definitions as  
769 Exemplified by the Definition of "Mammalia"." *Systematic Biology* 43: 497–510.

770

771 de Queiroz, K. 1997. "The Linnaean hierarchy and the evolutionization of taxonomy, with  
772 emphasis on the problem of nomenclature." *Aliso* 15:125–144.

773

774 de Queiroz, K., and J.A. Gauthier. 1990. "Phylogeny as a Central Principle in Taxonomy:  
775 Phylogenetic Definitions of Taxon Names." *Systematic Zoology* 39: 307–322.

776

777 de Queiroz, K., and J.A. Gauthier. 1992. "Phylogenetic Taxonomy." *Annual Review of Ecology*  
778 *and Systematics* 23: 449–480.  
779

780 de Queiroz K., and Gauthier J.A. 1994. "Toward a phylogenetic system of biological  
781 nomenclature." *Trends in Ecology and Evolution* 9: 27–31.  
782

783 de Queiroz, K., and M.J. Donoghue. 2011. "Phylogenetic Nomenclature, Three-Taxon  
784 Statements, and Unnecessary Name Changes." *Systematic Biology* 60: 887–892.  
785

786 de Queiroz, K., and M.J. Donoghue. 2013. "Phylogenetic Nomenclature, Hierarchical  
787 Information, and Testability." *Systematic Biology* 62: 167–174.  
788

789 De Queiroz, K., P.D. Cantino, and J.A. Gauthier. 2020. *Phylonyms: a companion to the*  
790 *PhyloCode*. Boca Raton, FL: CRC Press.  
791

792 Deans, A.R., M.J Yoder, and J.P. Balhoff. 2011. "Time to change how we describe biodiversity."  
793 *Trends in Ecology and Evolution* 27: 78–84.  
794

795 Deans, A.R., S.E. Lewis, E. Huala, S.S. Anzaldo, M. Ashburner, J.P. Balhoff, D.C. Blackburn, et  
796 al. 2015. "Finding Our Way through Phenotypes." *PLoS Biology* 13: e1002033.  
797

798 Eliason, C.M., M.J. Andersen, S.J. Hackett. 2019. "Using Historical Biogeography Models to  
799 Study Color Pattern Evolution." *Systematic Biology* 68: 755-766.  
800

801 Ereshefsky, M. 2001. *The Poverty of the Linnaean Hierarchy. A philosophical study of biological*  
802 *taxonomy*. Cambridge (MA): Cambridge University Press.  
803

804 Eriksson, T., M.J. Donoghue, and M.S. Hibbs. 1998. "Phylogenetic analysis of *Potentilla* using  
805 DNA sequences of nuclear ribosomal internal transcribed spacers (ITS), and implications for  
806 the classification of Rosoideae (Rosaceae)." *Plant Systematics and Evolution* 211: 155–179.  
807

808 Estes, R., K. de Queiroz, and J. Gauthier. 1988. "Phylogenetic relationships within Squamata."  
809 In *Phylogenetic relationships of the lizard families: essays commemorating Charles L.*  
810 *Camp*, edited by R. Estes, and G.K. Pregill, 119–281. Stanford, CA: Stanford University  
811 Press.

812

813 Folk, R.A., R.L. Stubbs, M.E. Mort, N. Cellinese, J.M. Allen, P.S. Soltis, D.E. Soltis, and R.P.  
814 Guralnick. 2019. "Rates of niche and phenotype evolution lag behind diversification in a  
815 temperate radiation." *Proceedings of the National Academy of Science* 116: 10874–10882.

816

817 Franz, N.M., M. Chen, P. Kianmajd, S. Yu, S. Bowers, A.S. Weakley, and B. Ludäscher. 2016.  
818 "Names are not good enough: reasoning over taxonomic change in the *Andropogon* complex  
819 1." *Semantic Web* 7: 645–667.

820

821 Franz, N.M., L.J. Musher, J.W. Brown, S. Yu, and B. Ludäscher. 2019. "Verbalizing  
822 phylogenomic conflict: Representation of node congruence across competing reconstructions  
823 of the neoavian explosion." *PLoS Computational Biology* 15: e1006493.

824

825 Gauthier, J., and K. Padian. 1985. "Phylogenetic, functional, and aerodynamic analyses of the  
826 origin of birds and their flight." In *The beginnings of birds*, edited by K. Hecht, G.H. Ostrom,  
827 G. Viohl, P. Wellnhofer, 185–197. Eichstatt (Germany): Freude des Jura-Museums.

828

829 Gauthier, J. 1986. "Saurischian monophyly and the origin of birds." In *The origin of birds and*  
830 *the evolution of flight*, edited by K. Padian, 1–55. San Francisco, CA: California Academy of  
831 Sciences.

832

833 Gauthier, J., R. Estes, and K. de Queiroz. 1988. "A phylogenetic analysis of  
834 Lepidosauromorpha." In *Phylogenetic relationships of the lizard families: essays*  
835 *commemorating Charles L. Camp*, edited by R. Estes, and G.K. Pregill, 15–98. Stanford, CA:  
836 Stanford University Press.

837

838 Ghiselin, M.T. 1984. "Definition, Character, and Other Equivocal Terms." *Systematic Zoology*  
839 33: 104–110.  
840

841 Haendel, M., J. Balhoff, F. Bastian, D. Blackburn, J. Blake, Y. Bradford, A. Comte, et al. 2014.  
842 "Unification of multi-species vertebrate anatomy ontologies for comparative biology in  
843 Uberon." *Journal of Biomedical Semantics* 5: 21.  
844

845 Hampton, S.E., M.B. Jones, L.W. Wasser, M.P. Schildhauer, S.R. Supp, J. Brun, R.R.  
846 Hernandez, et al. 2017. "Skills and Knowledge for Data-Intensive Environmental Research."  
847 *BioScience* 67: 546–557.  
848

849 Härlin, M., and P. Sundberg. 1998. "Taxonomy and Philosophy of Names." *Biology and*  
850 *Philosophy* 13: 233–244.  
851

852 Hennig, W. 1950. *Grundzüge einer Theorie der phylogenetischen Systematik*. Berlin: Deutscher  
853 Zentralverlag.  
854

855 Hennig, W. 1966. *Phylogenetic Systematics*. Urbana, IL: University of Illinois Press.

856 Hibbett, D.S., and M.J. Donoghue. 1998. "Integrating Phylogenetic Analysis and Classification  
857 in Fungi." *Mycologia* 90: 347–356.  
858

859 Hibbett, D.S., R.H. Nilsson, M. Snyder, M. Fonseca, J. Costanzo, M. Shonfeld. 2005.  
860 "Automated Phylogenetic Taxonomy: An Example in the Homobasidiomycetes (Mushroom-  
861 Forming Fungi)." *Systematic Biology* 54: 660–668.  
862

863 Hibbett, D.S. 2016. "The invisible dimension of fungal diversity." *Science* 351: 1150-1151.  
864

865 Hibbett, D.S., M. Blackwell, T.Y. James, J.W. Spatafora, J.W. Taylor, and R. Vilgalys. 2018.  
866 "Phylogenetic taxon definitions for Fungi, Dikarya, Ascomycota and Basidiomycota." *IMA*  
867 *Fungus* 9: 291–298.  
868

869 Hill M.S., A.L. Hill, J. Lopez, K.J. Peterson, S. Pomponi, M.C. Diaz, R.W. Thacker, et al. 2013.  
870 “Reconstruction of Family-Level Phylogenetic Relationships within Demospongiae  
871 (Porifera) Using Nuclear Encoded Housekeeping Genes.” *PLoS ONE* 8: e50437.  
872

873 Hinchliff C.E., S.A. Smith, J.F. Allman, J.G. Burleigh, R. Chaudhary, L.M. Coghill, K.A.  
874 Crandall, et al. 2015. “Synthesis of phylogeny and taxonomy into a comprehensive tree of  
875 life.” *Proceeding of the National Academy of Science* 112: 12764–12769.  
876

877 Horridge, M., and P.F. Patel-Schneider. 2012. OWL 2 Web Ontology Language Manchester  
878 Syntax (Second Edition). <https://www.w3.org/TR/owl2-manchester-syntax/>  
879

880 Hortal, J., F. de Bello, J.A.F. Diniz-Filho, T.M. Lewinsohn, J.M. Lobo, and R.J. Ladle. 2015.  
881 “Seven Shortfalls that Beset Large-Scale Knowledge of Biodiversity.” *Annual Review of*  
882 *Ecology, Evolution, and Systematics* 46: 523–549.  
883

884 Howard, C.C., R. Folk, J.M. Beaulieu, and N. Cellinese. 2019. “The monocotyledonous  
885 underground: global climatic and phylogenetic patterns of geophyte diversity.” *American*  
886 *Journal of Botany* 106: 850–863.

887 Howard, C.C., J.B. Landis, J.M. Beaulieu, and N. Cellinese. 2020. “Geophytism in monocots  
888 lead to higher rates of diversification.” *New Phytologist* 225: 1023-1032.  
889

890 Hundt, P.J., S.P. Iglésias, A.S. Hoey, and A.M. Simons. 2014. “A multilocus molecular  
891 phylogeny of combtooth blennies (Percomorpha: Blennioidei: Blenniidae): Multiple  
892 invasions of intertidal habitats.” *Molecular Phylogenetics and Evolution* 70: 47–56.  
893

894 Jensen, L.J., and P. Bork. 2010. “Ontologies in quantitative biology: a basis for comparison,  
895 integration, and discovery.” *PLoS Biology* 8: e1000374.  
896

897 Joyce, W.G., J.F. Parham, and J.A. Gauthier. 2004. “Developing a protocol for the conversion of  
898 rank-based taxon names to phylogenetically defined clade names, as exemplified by turtles.”  
899 *Journal of Paleontology* 78: 989–1013.

900  
901 Judd, W.S., W. Stern, V.I., and Cheadle. 1993. "Phylogenetic position of *Apostasia* and  
902 *Neuwiedia* (Orchidaceae)." *Botanical Journal of the Linnaean Society* 113: 87–94.  
903  
904 Judd, W.S., R.W. Sanders, and M.J. Donoghue. 1994. "Angiosperm family pairs: preliminary  
905 phylogenetic analyses." *Harvard Papers in Botany* 5: 1–51.  
906  
907 Keeseey, T.M. 2007. "A mathematical approach to defining clade names, with potential  
908 applications to computer storage and processing." *Zoologica Scripta* 36: 607–621.  
909  
910 Kim, O.-S., Y.-J. Cho, K. Lee, S.-H. Yoon, M. Kim, H. Na, S.-C. Park, et al.. 2012. "Introducing  
911 EzTaxon-e: a prokaryotic 16S rRNA gene sequence database with phylotypes that represent  
912 uncultured species." *International Journal of Systematics, Evolution and Microbiology* 62:  
913 716–721.  
914  
915 Kozlov, A.M., D. Darriba, T. Flouri, B. Morel, and A. Stamatakis. 2019. "RAxML-NG: A fast,  
916 scalable, and user-friendly tool for maximum likelihood phylogenetic inference."  
917 *Bioinformatics* 35: 4453-4455.  
918 Kron, K.A. 1997. "Exploring alternative systems of classification." *Aliso* 15: 105–111.  
919  
920 Lammers, T.G. 2007. "Campanulaceae." In *The Families and Genera of Vascular Plants*,  
921 edited by J.W. Kadereit, and C. Jeffrey, 8: 26–56. Berlin, Heidelberg: Springer Verlag.  
922  
923 Lee, M.S.Y. 1996. "Stability in Meaning and Content of Taxon Names: An Evaluation of  
924 Crown-Clade Definitions." *Proceedings of the Royal Society of London Series B*. 263: 1103–  
925 1109.  
926  
927 Lee, M.S.Y. 1998. "Phylogenetic Uncertainty, Molecular Sequences, and the Definition of  
928 Taxon Names." *Systematic Biology* 47: 719–726.  
929  
930 Lee, M.S.Y. 2001. "On Recent Arguments for Phylogenetic Nomenclature." *Taxon* 50: 175–180.

931  
932 Lemmon, E.M., and A.R. Lemmon. 2013. “High-Throughput Genomic Data in Systematics and  
933 Phylogenetics.” *Annual Review of Ecology, Evolution and Systematics* 44: 99–121.  
934  
935 Lin, C.H., K.C. Tsai, P. Prior, J.F. and Wang. 2014. “Phylogenetic relationships and population  
936 structure of *Ralstonia solanacearum* isolated from diverse origins in Taiwan.” *Plant*  
937 *Pathology* 63: 1395–1403.  
938  
939 Linnaeus, C. 1753. *Species Plantarum*. Stockholm: Laurentius Salvius.  
940  
941 Mabee, P., J.P. Balhoff, W.M. Dahdul, H. Lapp, P.E. Midford, T.J. Vision, and M. Westerfield.  
942 2012. “500,000 fish phenotypes: The new informatics landscape for evolutionary and  
943 developmental biology of the vertebrate skeleton.” *Journal of Applied Ichthyology* 28: 300–  
944 305.  
945  
946 Mabee, P.M., W.M. Dahdul, J.P. Balhoff, H. Lapp, P. Manda, J. Uyeda, T. Vision, and M.  
947 Westerfield. 2018. “Phenoscape: Semantic analysis of organismal traits and genes yields  
948 insights in evolutionary biology.” In *Application of Semantic Technology in Biodiversity*  
949 *Science*, edited by A.E. Thessen, 207-224. Berlin: IOS Press.  
950  
951 Madzia, D., and A. Cau. 2017. “Inferring ‘weak spots’ in phylogenetic trees: application to  
952 mosasauroid nomenclature.” *PeerJ* 5: e3782.  
953  
954 Manda, P., J.P. Balhoff, H. Lapp, P. Mabee, and T.J. Vision. 2015. “Using the Phenoscape  
955 Knowledgebase to relate genetic perturbations to phenotypic evolution.” *Genesis* 53: 561–  
956 571.  
957  
958 Mannion, P.D., P. Upchurch, R.N. Barnes, and O. Mateus. 2013. “Osteology of the Late Jurassic  
959 Portuguese sauropod dinosaur *Lusotitan atalaiensis* (Macronaria) and the evolutionary  
960 history of basal titanosauriforms.” *Zoological Journal of the Linnaean Society* 168: 98–206.  
961



962 Massana, R., E.F. DeLong, and C. Pedros-Alio. 2000. "A few cosmopolitan phylotypes dominate  
963 planktonic archaeal assemblages in widely different oceanic provinces." *Applied and*  
964 *Environmental Microbiology* 66: 1777–1787.

965

966 McTavish, E.J., B.T. Drew, B. Redelings, and K.A. Cranston. 2017. "How and Why to Build a  
967 Unified Tree of Life." *BioEssays* 39: 1700114.

968

969 Michener, W.K., and M.B. Jones. 2012. "Ecoinformatics: supporting ecology as a data-intensive  
970 science." *Trends in Ecology and Evolution* 27: 85–93.

971

972 Mishler, B., and J. S. Wilkins. 2018. "The Hunting of the SNaRC: A Snarky Solution to the  
973 Species Problem." *Philosophy, Theory, and Practice in Biology* 10.  
974 <https://doi.org/10.3998/ptpbio.16039257.0010.001>

975

976 Mungall, C.J., G.V. Gkoutos, C.L. Smith, M.A. Haendel, S.E. Lewis, and M. Ashburner. 2010.  
977 "Integrating phenotype ontologies across multiple species." *Genome Biology* 11: R2  
978 <https://doi.org/10.1186/gb-2010-11-1-r2>

979 Mungall, C.J., M. Bada, T.Z. Berardini, J. Deegan, A. Ireland, M.A. Harris, D.P. Hill, and J.  
980 Lomax. 2011. "Cross-product extensions of the Gene Ontology." *Journal of Biomedical*  
981 *Informatics* 44: 80–86.

982

983 Mungall, C.J., C. G.V. Torniai, Gkoutos, S.E. Lewis, and M.A. Haendel. 2012. "Uberon, an  
984 integrative multi-species anatomy ontology." *Genome Biology* 13: R5  
985 <https://doi.org/10.1186/gb-2012-13-1-r5>

986

987 Page, L.M., B.J. MacFadden, J.A. Fortes, P.S. Soltis, and G. Riccardi. 2015. "Digitization of  
988 Biodiversity Collections Reveals Biggest Data on Biodiversity." *BioScience* 65: 841-842.

989

990 Parr, C.S., R. Guralnick, N. Cellinese, and R.D.M. Page. 2012. "Evolutionary informatics:  
991 unifying knowledge about the diversity of life." *Trends in Ecology and Evolution* 27: 94–  
992 103.

993  
994 Pesquita, C., D. Faria, A.O. Falcão, P. Lord, and F.M. Couto. 2009. “Semantic similarity in  
995 biomedical ontologies.” *PLoS Computational Biology* 5: e1000443.  
996  
997 Philippe, H., F. Delsuc, H. Brinkmann, and N. Lartillot. 2005. “Phylogenomics.” *Annual Review*  
998 *of Ecology, Evolution and Systematics* 36: 541–562.  
999  
1000 Pleijel, F. 1999. “Phylogenetic Taxonomy, a Farewell to Species, and a Revision of  
1001 Heteropodarke (Hesionidae, Polychaeta, Annelida).” *Systematic Biology* 48: 755–789.  
1002  
1003 Polaszek, A., and E.O. Wilson. 2005. “Sense and stability in animal names.” *Trends in Ecology*  
1004 *and Evolution* 20: 421–422.  
1005  
1006 Porter, J.H., E. Nagy, T.K. Kratz, P. Hanson, S.L. Collins, and P. Arzberger. 2009. “New eyes on  
1007 the world: Advanced sensors for ecology.” *BioScience* 59: 385–397.  
1008  
1009 Prosdocimi, F., B. Chisham, E. Pontelli, J.D. Thompson, and A. Stoltzfus. 2009. “Initial  
1010 Implementation of a comparative Data Analysis Ontology.” *Evolutionary Bioinformatics*  
1011 doi:[10.4137/EBO.S2320](https://doi.org/10.4137/EBO.S2320)  
1012  
1013 Rabi, M., V.B. Sukhanov, V.N. Egorova, I. Danilov, and W.G. Joyce. 2014. “Osteology,  
1014 relationships, and ecology of *Annemys* (Testudines, Eucryptodira) from the Late Jurassic of  
1015 Shar Teg, Mongolia, and phylogenetic definitions for Xinjiangchelyidae, Sinemydidae, and  
1016 Macrobaenidae.” *Journal of Vertebrate Paleontology* 34: 327–352.  
1017  
1018 Randell, D. A., Z. Cui, and A. Cohn. 1992. “A spatial logic based on regions and connection.” In  
1019 *Principles of Knowledge Representation and Reasoning: Proceedings of the Third*  
1020 *International Conference*, edited by B. Nebel, C. Rich, and W. Swartout, 165–176. San  
1021 Mateo, CA: Morgan Kaufmann.  
1022

1023 Rees, J.A., and K. Cranston. 2017. "Automated assembly of a reference taxonomy for  
1024 phylogenetic data synthesis." *Biodiversity Data Journal* e12581 doi:[10.3897/BDJ.5.e12581](https://doi.org/10.3897/BDJ.5.e12581)  
1025

1026 Rieppel, O. 2006. "The PhyloCode: a critical discussion of its theoretical foundation." *Cladistics*  
1027 22: 186–197.  
1028

1029 Rowe, T. 1987. "Definition and Diagnosis in the Phylogenetic System." *Systematic Zoology* 36:  
1030 208–211.  
1031

1032 Rowe, T., and J. Gauthier. 1992. "Ancestry, paleontology and definition of the name  
1033 Mammalia." *Systematic Biology* 41: 372–378.  
1034

1035 Schander, C., and M. Thollessen. 1995. "Phylogenetic taxonomy-some comments." *Zoologica*  
1036 *Scripta* 24: 263–268.  
1037

1038 Senderov, V., K. Simov, N. Franz, P. Stoev, T. Catapano, D. Agosti, G. Sautter, R.A. Morris,  
1039 and L. Penev. 2018. "OpenBiodiv-O: ontology of the OpenBiodiv knowledge management  
1040 system." *Journal of Biomedical Semantics* 9: 5 <https://doi.org/10.1186/s13326-017-0174-5>  
1041

1042 Schoch, R.R. 2013. "The evolution of major temnospondyl clades: an inclusive phylogenetic  
1043 analysis." *Journal of Systematic Palaeontology* 11: 673–705.  
1044

1045 Schuh, R.T. 2003. "The Linnaean system and its 250-year persistence." *Botanical Review* 69:  
1046 59–78.  
1047

1048 Sereno, P.C. 1999. "Definitions in Phylogenetic Taxonomy: Critique and Rationale." *Systematic*  
1049 *Biology* 48: 329–351.  
1050

1051 Sereno, P.C. 2005. "The Logical Basis of Phylogenetic Taxonomy." *Systematic Biology* 54:  
1052 595–619.  
1053

1054 Sereno, P.C., S. McAllister, and S.L. Brusatte. 2005. "TaxonSearch: a relational database for  
1055 suprageneric taxa and phylogenetic definitions." *Phyloinformatics* 8: 1–21.  
1056

1057 Sferco, E., A. López-Arbarello, and A.M. Báez. 2015. "Phylogenetic relationships of †*Luisiella*  
1058 *feruglioi* (Bordas) and the recognition of a new clade of freshwater teleosts from the Jurassic  
1059 of Gondwana." *BMC Evolutionary Biology* 15: 268.  
1060

1061 Smith, B., M. Ashburner, C. Rosse, J. Bard, W. Bug, W. Ceusters, L. Goldberg, et al. 2007. "The  
1062 OBO Foundry: coordinated evolution of ontologies to support biomedical data integration."  
1063 *Nature Biotechnology* 25: 1251–1255.  
1064

1065 Smith, S.A., J.M. Beaulieu, A. Stamatakis, and M.J. Donoghue. 2011. "Understanding  
1066 angiosperm diversification using small and large phylogenetic trees." *American Journal of*  
1067 *Botany* 98: 404–414.  
1068

1069 Smith, S.A. and J.W. Brown 2018. "Constructing a broadly inclusive seed plant phylogeny."  
1070 *American Journal of Botany* 105: 302-314.  
1071

1072 Soltis, D.E., S.A. Smith, N. Cellinese, K.J. Wurdack, D.C. Tank, S.F. Brockington, N.F. Refulio-  
1073 Rodriguez, et al. 2011. "Angiosperm phylogeny: 17 genes, 640 taxa." *American Journal of*  
1074 *Botany* 98: 704–730.  
1075

1076 Spatafora, J.W., Y. Chang, G.L. Benny, K. Lazarus, M.E. Smith, M.L. Berbee, G. Bonito, et al.  
1077 2016. "A phylum-level phylogenetic classification of zygomycete fungi based on genome-  
1078 scale data." *Mycologia* 108: 1028–1046.

1079 Sterli, J., D. Pol, and M. Laurin. 2013. "Incorporating phylogenetic uncertainty on phylogeny-  
1080 based palaeontological dating and the timing of turtle diversification." *Cladistics* 29: 233-  
1081 246.  
1082

1083 Sterner, B., and N.M. Franz. 2017. "Taxonomy for Humans or Computers? Cognitive Pragmatics  
1084 for Big Data." *Biological Theory* 12: 99–111.

1085  
1086 Stevens, P.F. 2006. “An end to all things?—plants and their names.” *Australian Systematic*  
1087 *Botany* 19: 115–133.  
1088  
1089 Stevens, P. F. 2017. “Angiosperm Phylogeny Website.” Version 14, July 2017  
1090 <http://www.mobot.org/MOBOT/research/APweb/>  
1091  
1092 Sundberg, P., and F. Pleijel. 1994. “Phylogenetic classification and the definition of taxon  
1093 names.” *Zoologica Scripta* 23: 19–25.  
1094  
1095 Thau, D., and B. Ludäscher 2007. “Reasoning about taxonomies in first-order logic.” *Ecological*  
1096 *informatics* 2: 195–209.  
1097  
1098 Thau, D., S. Bowers, and B. Ludäscher 2008. “Merging taxonomies under RCC-5 algebraic  
1099 articulations.” *Proceedings of the 2nd international workshop on Ontologies and information*  
1100 *systems for the semantic web*. 47–54 doi:10.1145/1458484.1458492  
1101  
1102 Thessen, A.E., D.E. Bunker, P.L. Buttigieg, L.D. Cooper, W.M. Dahdul, S. Domisch, N.M.  
1103 Franz, et al. 2015. “Emerging semantics to link phenotype and environment.” *PeerJ*. 3:  
1104 e1470 doi: [10.7717/peerj.1470](https://doi.org/10.7717/peerj.1470)  
1105  
1106 Torres-Carvajal, O., and P. Mafla-Endara. 2013. “Evolutionary history of Andean Pholidobolus  
1107 and Macropholidus (Squamata: Gymnophthalmidae) lizards.” *Molecular Phylogenetics and*  
1108 *Evolution* 68: 212–217.  
1109  
1110 Vision, T., J. Blake, H. Lapp, P. Mabee, and M. Westerfield. 2011. “Similarity between semantic  
1111 description sets: addressing needs beyond data integration.” In *Proceedings of the First*  
1112 *International Workshop on Linked Science (LISC 2011)*, edited by T. Kauppinen, L.C.  
1113 Pouchard, and C. Keßler. Bonn: CEUR Workshop Proceedings.  
1114

1115 Vogt, L. 2009. "The future role of bio-ontologies for developing a general data standard in  
1116 biology: chance and challenge for zoo-morphology." *Zoomorphology* 128: 201–217.  
1117

1118 W3C OWL Working Group. 2012. *OWL 2 Web Ontology Language Document Overview*  
1119 *(second edition)*. <https://www.w3.org/TR/owl2-overview>.  
1120

1121 Washington, N.L., M.A. Haendel, C.J. Mungall, M. Ashburner, M. Westerfield, and S.E. Lewis.  
1122 2009. "Linking human diseases to animal models using ontology-based phenotype  
1123 annotation." *PLoS Biology* 7: e1000247 doi:10.1371/journal.pbio.1000247  
1124

1125 Wojciechowski, M.F. 2013. "Towards a new classification of Leguminosae: Naming clades  
1126 using non-Linnaean phylogenetic nomenclature." *South African Journal of Botany* 89: 85–  
1127 93.  
1128

1129 Wright, D.F., W.I. Ausich, S.R. Cole, M.E. Peter, and E.C. Rhenberg. 2017. "Phylogenetic  
1130 taxonomy and classification of the Crinoidea (Echinodermata)." *Journal of Paleontology* 91:  
1131 829–846.  
1132

1133 Wyss, A.R., and J. Meng. 1996. "Application of phylogenetic taxonomy to poorly resolved  
1134 crown clades: a stem-modified node-based definition of Rodentia." *Systematic Biology* 45:  
1135 559–568.  
1136

1137 Zimmermann, W. 1931 (1937). "Arbeitsweise der botanischen Phylogenetik und anderer  
1138 Gruppierungswissenschaften." In *Handbuch der biologischen Arbeitsmethoden*, edited by E.  
1139 Abderhalden, 941–1053. Berlin: Urban & Schwarzenberg.  
1140

1141 Zimmermann, W. 1934. "Research on phylogeny of species and of single characters." *American*  
1142 *Naturalist* 68: 381–384.  
1143

1144 Zimmermann, W. 1943. "Die Methoden der Phylogenetik." In *Evolution der Organismen*, edited  
1145 by G. Heberer, 20–56. Jena: Fischer.

1146

1147

1148 **Figure captions**

1149

1150 Figure 1. Upper half: phylogeny of Campanuloideae redrawn from Crowl et al. (2016) showing  
1151 the polyphyly of *Campanula* (lineages in blue). Lower half: Distribution of *Campanula* as  
1152 retrieved from a GBIF query.

1153

1154 Figure 2. The three basic clade definitions.

1155

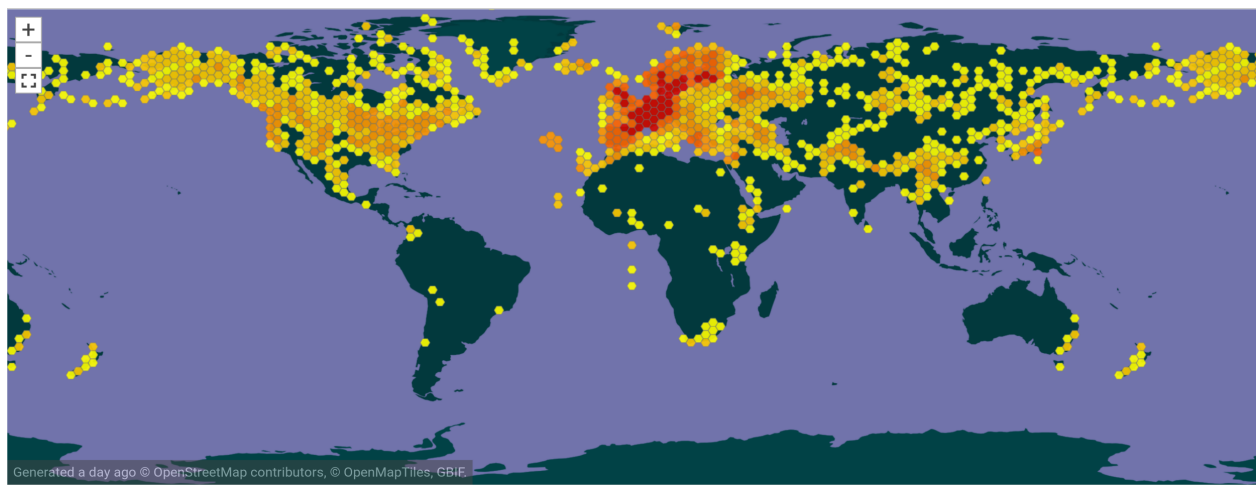
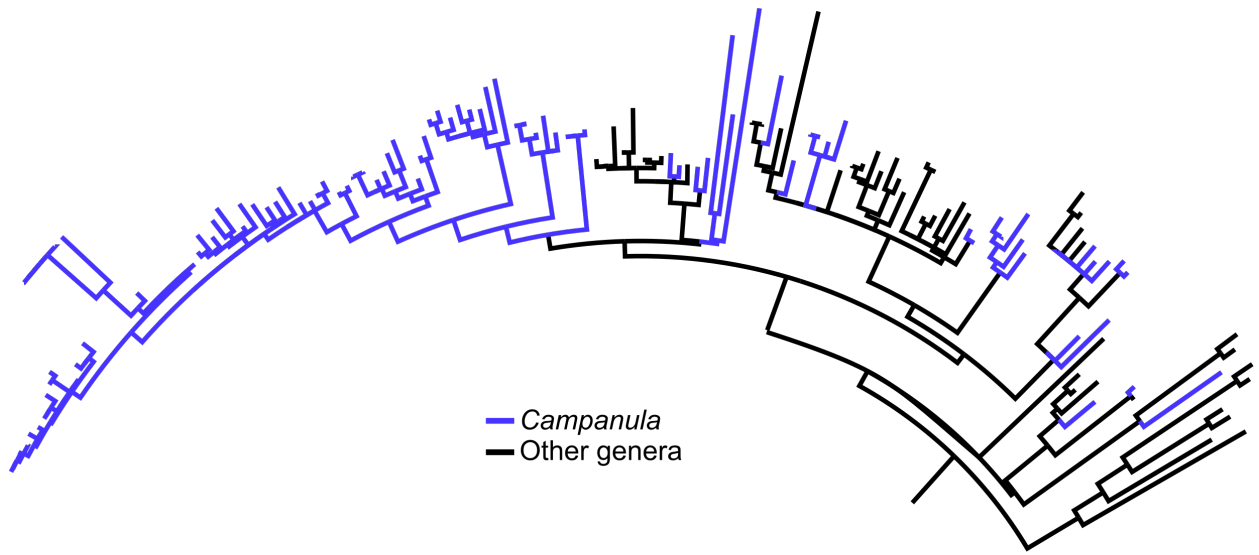
1156 Figure 3. Phylogeny of Asterales showing the clade Campanulaceae with its five lineages, the  
1157 sister group Rouseaceae, and other related lineages (adapted from Steven 2017).

1158

1159

1160 Figure 1.

1161

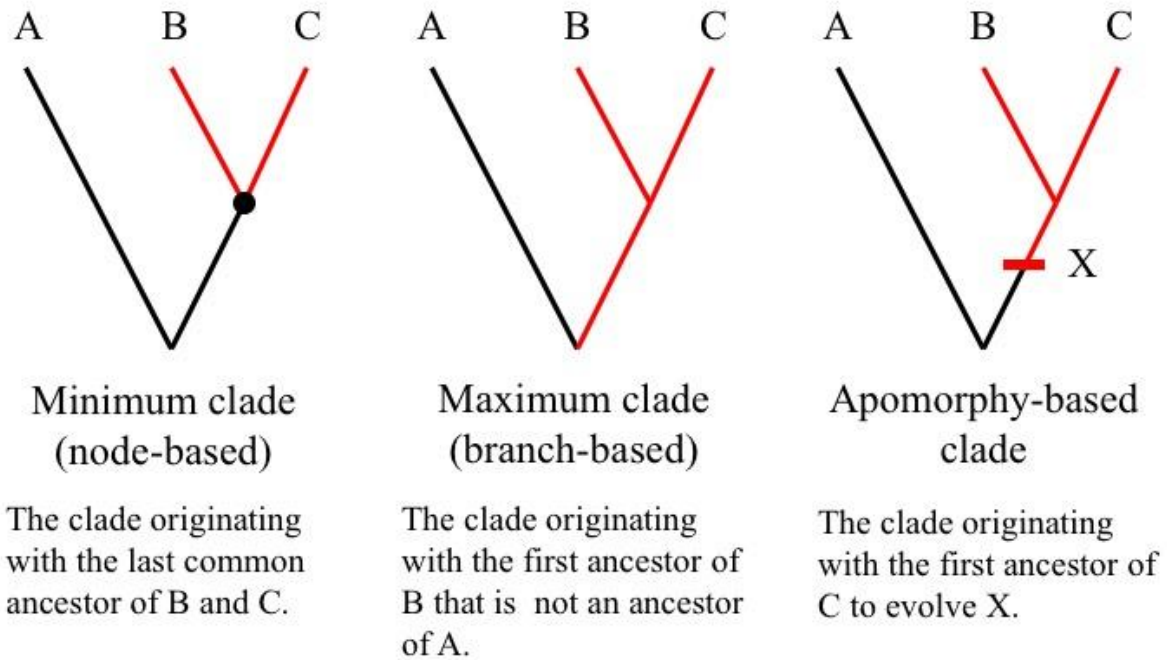


1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174

Figure 2.

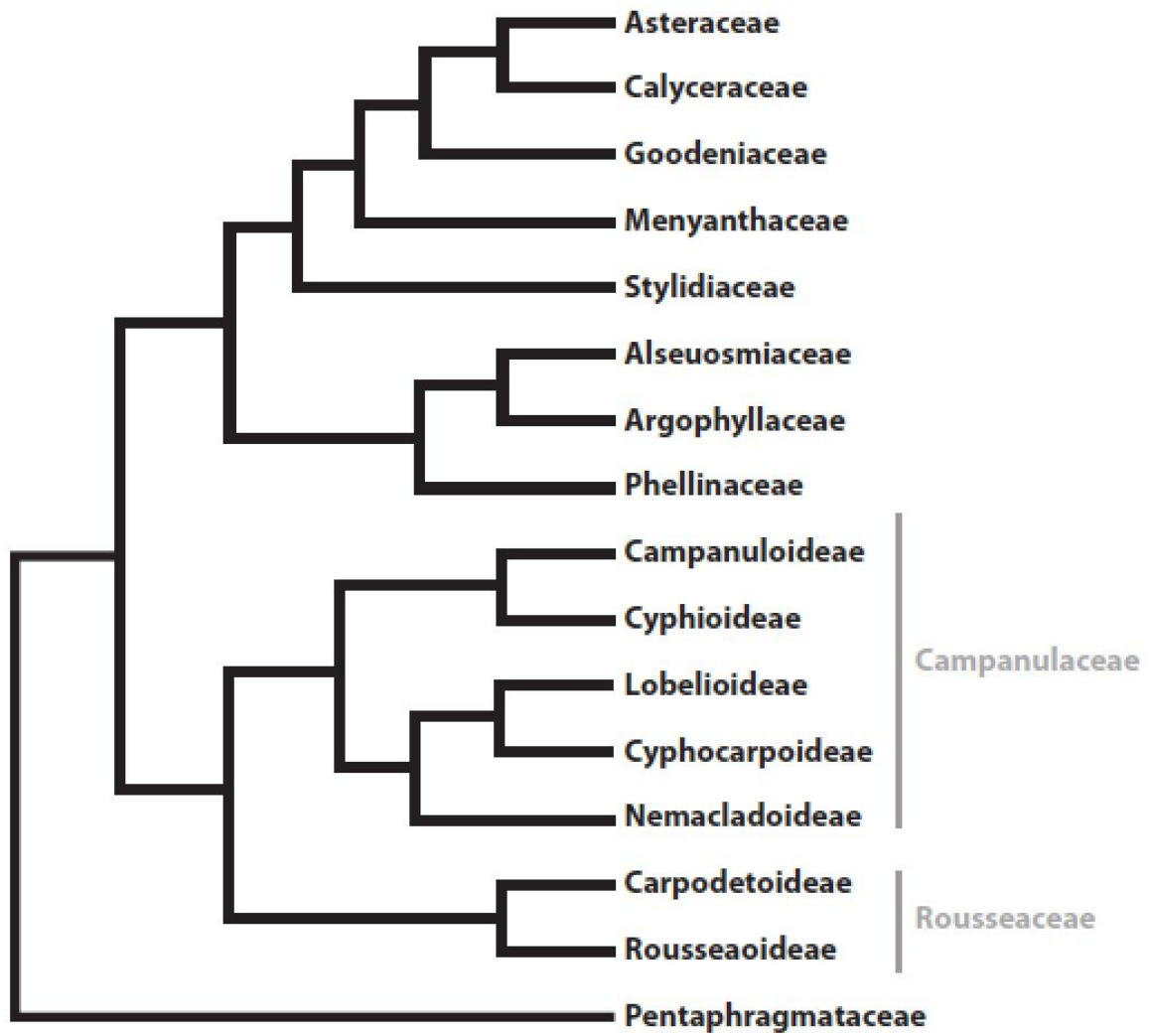


# Phylogenetic Definitions



1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189

Figure 3.



1190