

# RNeXML: a package for reading and writing richly annotated phylogenetic, character, and trait data in R

Carl Boettiger<sup>\*,a</sup>, Scott Chamberlain<sup>c</sup>, Rutger Vos<sup>d</sup>, Hilmar Lapp<sup>b</sup>

<sup>a</sup>*Center for Stock Assessment Research, Department of Applied Math and Statistics, University of California, Mail Stop SOE-2, Santa Cruz, CA 95064, USA*

<sup>b</sup>*National Evolutionary Synthesis Center and Duke University, Durham, NC, USA*

<sup>c</sup>*University of California, Berkeley, CA, USA*

<sup>d</sup>*Naturalis Biodiversity Center, Leiden, the Netherlands*

## Abstract

1. NeXML is a powerful and forward-compatible exchange standard recently proposed to better meet the expanding needs for machine-proof phylogenetic data and metadata sharing. Here we present the RNeXML package, which aims to provide to users of the popular R programming language easy programmatic access to reading and writing NeXML documents, including rich metadata, in a way that interfaces seamlessly with the extensive library of phylogenetic tools already available in the R ecosystem.
2. Wherever possible, we designed RNeXML to map NeXML document contents, whose arrangement is influenced by the format's normative XML Schema definition, to their most intuitive or useful representation in R. To make NeXML's powerful facility for recording semantically rich machine-readable metadata accessible to R users, we designed a functional programming interface to it that hides the semantic web standards leveraged by NeXML from R users who are unfamiliar with them.
3. RNeXML can read any NeXML document that validates, and it generates valid NeXML documents from phylogeny and character data in various R representations in use. The metadata programming interface at a basic level aids fulfilling data documentation best practices, and at an advanced level preserves NeXML's nearly limitless extensibility, for which we provide a fully working demonstration. Furthermore, to lower the barriers to sharing well-documented phylogenetic data, RNeXML has started to integrate with taxonomic metadata augmentation services on the web, and with online repositories for data archiving.
4. RNeXML allows R's rich ecosystem to read and write data in the NeXML format through an interface that is no more involved than reading or writing data from other, less powerful data

---

\*Corresponding author  
Email address: `cboettig(at)gmail.com` (Carl Boettiger)

formats. It also provides an interface designed to feel familiar to R programmers and to be consistent with recommended practices for R package development, yet that retains the full power for users to add their own custom data and metadata to the phylogenies they work with, without introducing potentially incompatible changes to the exchange standard.

## Introduction

Users of the popular statistical and mathematical computing platform R (R Core Team 2014) enjoy a wealth of readily installable comparative phylogenetic methods and tools (O’Meara 2014). Exploiting the opportunities arising from this wealth for complex and integrative comparative research questions relies on the ability to reuse and integrate previously generated or published data and metadata. The expanding data exchange needs of the evolutionary research community are rapidly outpacing the capabilities of most current and widely used data exchange standards (Vos *et al.* 2012), which were all developed a decade or more ago. This has resulted in a radiation of different data representations and exchange standard “flavors” that are no longer interoperable at the very time when the growth of available data and methods has made that interoperability most valuable. In response to the unmet needs for standardized data exchange in phylogenetics, a modern XML-based exchange standard, called NeXML, has recently been developed (Vos *et al.* 2012). NeXML comprehensively supports current data exchange needs, is predictably machine-readable, and is forward compatible.

The exchange problem for phylogenetic data is particularly acute in light of the challenges in finding and sharing phylogenetic data without the otherwise common loss of most data and metadata semantics (Stoltzfus *et al.* 2012; Drew *et al.* 2013; Cranston *et al.* 2014). For example, the still popular NEXUS file format (Maddison *et al.* 1997) cannot consistently represent horizontal gene transfer or ambiguity in reading a character (such as a DNA sequence base pair). This and other limitations have led to modifications of NEXUS in different ways for different needs, with the unfortunate result that NEXUS files generated by one program can be incompatible with another (Vos *et al.* 2012). Without a formal grammar, software based on NEXUS files may also make inconsistent assumptions about tokens, quoting, or element lengths. Vos *et al.* (2012) estimates that as many as 15% of the NEXUS files in the CIPRES portal contain unrecoverable but hard to diagnose errors.

A detailed account of how the NeXML standard addresses these and other relevant challenges can be found in Vos *et al.* (2012). In brief, NeXML was designed with the following important properties. First,

NeXML is defined by a precise grammar that can be programmatically **validated**. This ensures that any NeXML file can be expected to be parsed and read without errors by any NeXML-compliant software, provided the file passes validation. Second, NeXML is **extensible**: a user can define representations of new, previously unanticipated information (as we will illustrate) without violating its defining grammar. Third and most importantly, NeXML is **rich in computable semantics**. The formal semantics used in NeXML is machine-readable in a predictable manner, and also enables machine reasoning about the data and metadata (see Parr et al. (2011)), especially when metadata attributes and values are anchored in appropriately axiomatized ontologies.

To make the capabilities of NeXML available to R users in an easy-to-use form, and to lower the hurdles to adoption of the standard, we present RNeXML, an R package that aims to provide easy programmatic access to reading and writing NeXML documents, tailored for the kinds of use-cases that will be common for users and developers of the wealth of evolutionary analysis methods within the R ecosystem.

## The RNeXML package

The RNeXML package is written entirely in R and available under a Creative Commons Zero public domain waiver. The current development version can be found on Github at <https://github.com/ropensci/RNeXML>, and the stable version can be installed from the CRAN repository. RNeXML is part of the rOpenSci project. Users of RNeXML are encouraged to submit bug reports or feature requests in the issues log on Github, or the phylogenetics R users group list at [r-sig-phylo@r-project.org](mailto:r-sig-phylo@r-project.org) for help. Vignettes with more detailed examples of specific features of RNeXML are distributed with the R package and serve as a supplement to this manuscript. Each of the vignettes can be found at <http://ropensci.github.io/RNeXML/>.

## *Representation of NeXML documents in R*

Conceptually, a NeXML document has the following components: (1) phylogeny topology data, (2) character or trait data in matrix form, (3) operational taxonomic units (OTUs), and (4) metadata. To represent the contents of a NeXML document (currently in memory), RNeXML defines the `nexml` object type. This type therefore holds phylogenetic trees as well as character or trait matrices, and all metadata, which is similar to the phylogenetic data object types defined in the `phylobase` package

88 (NESCENT R Hackathon Team 2014), but contrasts with the more widely used ones defined in the `ape`  
89 package (Paradis *et al.* 2004), which represents trees alone.

90 When reading and writing NeXML documents, `RNeXML` aims to map their components to and from,  
91 respectively, their most widely used representations in R. As a result, the types of objects accepted  
92 or returned by the package's methods are the `phylo` and `multiPhylo` objects from the `ape` package  
93 (Paradis *et al.* 2004) for phylogenies, and R's native `data.frame` list structure for data matrices.

#### 94 *Reading phylogenies and character data*

95 The method `nexml_read()` reads NeXML files, either from a local file, or from a remote location via  
96 its URL, and returns an object of type `nexml`:

```
nex <- nexml_read("components/trees.xml")
```

97 The method `get_trees_list()` can be used to extract the phylogenies as an `ape::multiPhylo` object,  
98 which can be treated as a list of `ape::phylo` objects:

```
phy <- get_trees_list(nex)
```

99 The `get_trees_list()` method is designed for use in scripts, providing a consistent and predictable  
100 return type regardless of the number of phylogenies a NeXML document contains. For greater convenience  
101 in interactive use, the method `get_trees()` returns the R object most intuitive given the arrangement  
102 of phylogeny data in the source NeXML document. For example, the method returns an `ape::phylo`  
103 object if the NeXML document contains a single phylogeny, an `ape::multiPhylo` object if it contains  
104 multiple phylogenies arranged in a single `trees` block, and a list of `ape::multiPhylo` objects if it  
105 contains multiple `trees` blocks (the capability for which NeXML inherits from NEXUS).

106 If the location parameter with which the `nexml_read()` method is invoked is recognized as a URL,  
107 the method will automatically download the document to the local working directory and read it from  
108 there. This gives convenient and rapid access to phylogenetic data published in NeXML format on the  
109 web, such as the content of the phylogenetic data repository TreeBASE (Piel *et al.* 2002). For example,  
110 the following plots a tree in TreeBASE (using `ape`'s `plot` function):

```
tb_nex <- nexml_read(
  "https://raw.githubusercontent.com/TreeBASE/supertreebase/master/data/treebase/S100.xml")
tb_phy <- get_trees_list(tb_nex)
plot(tb_phy[[1]])
```

111 The method `get_characters()` obtains character data matrices from a `nexml` object, and returns them  
 112 as a standard `data.frame` R object with columns as characters and rows as taxa:

```
nex <- nexml_read("components/comp_analysis.xml")
get_characters(nex)
```

	log snout-vent length	reef-dwelling
113		
114	taxon_8	-3.2777799 0
115	taxon_9	2.0959433 1
116	taxon_10	3.1373971 0
117	taxon_1	4.7532824 1
118	taxon_2	-2.7624146 0
119	taxon_3	2.1049413 0
120	taxon_4	-4.9504770 0
121	taxon_5	1.2714718 1
122	taxon_6	6.2593966 1
123	taxon_7	0.9099634 1

124 A NeXML data matrix can be of molecular (for molecular sequence alignments), discrete (for most  
 125 morphological character data), or continuous type (for many trait data). To enable strict validation  
 126 of data types NeXML allows, and if their data types differ requires multiple data matrices to be  
 127 separated into different “blocks”. Since the `data.frame` data structure in R has no such constraints, the  
 128 `get_characters()` method combines such blocks as separate columns into a single `data.frame` object,  
 129 provided they correspond to the same taxa. Otherwise, a list of `data.frames` is returned, with list  
 130 elements corresponding to characters blocks. Similar to the methods for obtaining trees, there is also a  
 131 method `get_characters_list()`, which always returns a list of `data.frames`, one for each character  
 132 block.

### 133 *Writing phylogenies and character data*

134 The method `nexml_write()` generates a NeXML file from its input parameters. In its simplest  
135 invocation, the method writes a tree to a file:

```
data(bird.orders)
nexml_write(bird.orders, file = "birds.xml")
```

136 The first argument to `nexml_write()` is either an object of type `nexml`, or any object that can be coerced  
137 to it, such as in the above example an `ape::phylo` phylogeny. Alternatively, passing a `multiPhylo`  
138 object would write a list of phylogenies to the file.

139 In addition to trees, the `nexml_write()` method also allows to specify character data as another  
140 parameter. The following example uses data from the comparative phylogenetics R package `geiger`  
141 (Harmon *et al.* 2008).

```
library("geiger")
data(geospiza)
nexml_write(trees = geospiza$phy,
            characters = geospiza$dat,
            file="geospiza.xml")
```

### 142 *Validating NeXML*

143 File validation is a central feature of the NeXML format which ensures that any properly implemented  
144 NeXML parser will always be able to read the NeXML file. The function takes the path to any NeXML  
145 file and returns `TRUE` to indicate a valid file, or `FALSE` otherwise, along with a display of any error  
146 messages generated by the validator.

```
nexml_validate("geospiza.xml")
```

147 [1] TRUE

148 The `nexml_validate()` function performs this validation using the online NeXML validator (when a  
149 network connection is available), which performs additional checks not expressed in the NeXML schema  
150 itself (Vos *et al.* 2012). If a network connection is not available, the function falls back on the schema  
151 validation method from the `XML` package (Lang 2013).

## 152 *Creating and populating nexml objects*

153 Instead of packaging the various components for a NeXML file at the time of writing the file, RNeXML  
154 also allows users to create and iteratively populate in-memory nexml objects. The methods to do this  
155 are `add_characters()`, `add_trees()`, and `add_meta()`, for adding characters, trees, and metadata,  
156 respectively. Each of these functions will automatically create a new nexml object if not supplied with  
157 an existing one as the last (optional) argument.

158 For example, here we use `add_trees()` to first create a nexml object with the phylogeny data, and  
159 then add the character data to it:

```
nexObj <- add_trees(geospiza$phy)
nexObj <- add_characters(geospiza$dat, nexObj)
```

160 The data with which a nexml object is populated need not share the same OTUs. RNeXML automatically  
161 adds new, separate OTU blocks into the NeXML file for each data matrix and tree that uses a different  
162 set of OTUs.

163 Other than storage size, there is no limit to the number of phylogenies and character matrices that  
164 can be included in a single NeXML document. This allows, for example, to capture samples from a  
165 posterior probability distribution of inferred or simulated phylogenies and character states in a single  
166 NeXML file.

## 167 *Data documentation and annotation with built-in metadata*

168 NeXML allows attaching (“*annotating*”) metadata to any data element, and even to metadata themselves.  
169 Whether at the level of the document as a whole or an individual data matrix or phylogeny, metadata  
170 can provide bibliographic and provenance information, for example about the study as part of which  
171 the phylogeny was generated or applied, which data matrix and which methods were used to generate it.  
172 Metadata can also be attached to very specific elements of the data, such as specific traits, individual  
173 OTUs, nodes, or even edges of the phylogeny.

174 As described in Vos et al. (2012), to encode metadata annotations NeXML uses the “Resource  
175 Description Framework in Annotations” (RDFa) (Prud’hommeaux 2014). This standard provides for  
176 a strict machine-readable format yet enables future backwards compatibility with compliant NeXML

177 parsers (and thus RNeXML), because the capacity of a tool to *parse* annotations is not predicated on  
178 *understanding* the meaning of annotations it has not seen before.

179 To lower the barriers to sharing well-documented phylogenetic data, RNeXML aims to make recording  
180 useful and machine-readable metadata easier at several levels.

181 First, when writing a NeXML file the package adds certain basic metadata automatically if they  
182 are absent, using default values consistent with recommended best practices (Cranston *et al.* 2014).  
183 Currently, this includes naming the software generating the NeXML, a time-stamp of when a tree was  
184 produced, and an open data license.

185 Second, RNeXML provides a simple method, called `add_basic_metadata()`, to set metadata attributes  
186 commonly recommended for inclusion with data to be publicly archived or shared (Cranston *et al.*  
187 2014). The currently accepted parameters include `title`, `description`, `creator`, `pubdate`, `rights`,  
188 `publisher`, and `citation`. Behind the scenes the method automatically anchors these attributes in  
189 common vocabularies (such as Dublin Core).

190 Third, RNeXML integrates with the R package `taxize` (Chamberlain & Szöcs 2013) to mitigate one of  
191 the most common obstacles to reuse of phylogenetic data, namely the misspellings and inconsistent  
192 taxonomic naming with which OTU labels are often fraught. The `taxize_nexml()` method in RNeXML  
193 uses `taxize` to match OTU labels against the NCBI database, and, where a unique match is found, it  
194 annotates the respective OTU with the matching NCBI identifier.

#### 195 *Data annotation with custom metadata*

196 The RDFa standard adopted by NeXML (see Vos et al. (2012)) allows for arbitrary custom extensions to  
197 meet ongoing scientific needs without requiring modifications (or “extensions”) of the syntax. Specifically,  
198 every metadata annotation consists of a metadata property and a value of the property; the bearer  
199 of the property is determined by where the annotation appears in the NeXML document (technically,  
200 the containing element in the XML document model). This structure corresponds directly to Resource  
201 Description Framework (RDF) subject-predicate-object triples, and therefore the only requirement for  
202 custom metadata is that they be cast as RDF triples.

203 The RNeXML interface described above for built-in metadata annotations perform this cast automatically  
204 behind the scenes, including mapping metadata attributes to terms in requisite common vocabularies  
205 (such as Dublin Core for “title”, “creator”, etc.) or ontologies. Giving metadata properties (and



where applicable, values) as vocabulary or ontology terms rather than simple text strings is crucial for allowing machines to not only parse but also interpret and potentially reason over their semantics. To achieve this benefit for custom metadata extensions, the user necessarily needs to handle certain technical details from which the RNeXML API shields her otherwise, in particular the globally unique identifiers (normally HTTP URIs) of metadata terms and vocabularies. To be consistent with XML terminology, RNeXML calls vocabulary URIs *namespaces*, and their abbreviations *prefixes*. For example, the namespace for the Dublin Core Metadata Terms vocabulary is “<http://purl.org/dc/elements/1.1/>”. Using its common abbreviation “dc”, a metadata property “dc:title” expands to the identifier “<http://purl.org/dc/elements/1.1/title>”. This URI resolves to a human and machine-readable (depending on access) definition of precisely what the term `title` in Dublin Core means. In contrast, just using the text string “title” could also mean the title of a person, a legal title, the verb title, etc. URI identifiers of metadata vocabularies and terms are not mandated to resolve, but if machines are to derive the maximum benefit from them, they should resolve to a definition of their semantics in RDF.

RNeXML includes methods to obtain and manipulate metadata properties, values, identifiers, and namespaces. The `get_namespaces()` method accepts a `nexml` object and returns a named list of namespace prefixes and their corresponding identifiers known to the object:

```
birds <- nexml_read("birds.xml")
prefixes <- get_namespaces(birds)
prefixes["dc"]
```

```
                dc
"http://purl.org/dc/elements/1.1/"
```

The `get_metadata()` method returns, as a named list, the metadata annotations for a given `nexml` object at a given level, with the whole NeXML document being the default level (“all” extracts all metadata objects):

```
meta <- get_metadata(birds)
otu_meta <- get_metadata(birds, level="otu")
```

The returned list does not include the data elements to which the metadata are attached. Therefore, a different approach, documented in the metadata vignette, is recommended for accessing the metadata attached to data elements.

230 The `meta()` method creates a new metadata object from a property name and content (value). For  
231 example, the following creates a modification date metadata object, using a property in the PRISM  
232 vocabulary:

```
modified <- meta(property = "prism:modificationDate", content = "2013-10-04")
```

233 Metadata annotations in NeXML can be nested within another annotation, which the `meta()` method  
234 accommodates by accepting a parameter `children`, with the list of nested metadata objects (which can  
235 themselves be nested) as value.

236 The `add_meta()` function adds metadata objects as annotations to a `nexml` object at a specified level,  
237 with the default level being the NeXML document as a whole:

```
birds <- add_meta(modified, birds)
```

238 If the prefix used by the metadata property is not among the built-in ones (which can be obtained  
239 using `get_namespaces()`), it has to be provided along with its URI as the `namespaces` parameter. For  
240 example, the following uses the SKOS vocabulary (which currently is not built-in) to add a note to the  
241 trees in the `nexml` object:

```
history <- meta(property = "skos:historyNote",  
  content = "Mapped from the bird.orders data in the ape package using RNeXML")  
birds <- add_meta(history,  
  birds,  
  level = "trees",  
  namespaces = c(skos = "http://www.w3.org/2004/02/skos/core#"))
```

242 Alternatively, additional namespaces can also be added in batch using the `add_namespaces()` method.  
243 By virtue of subsetting the S4 `nexml` object, RNeXML also offers fine control of where a `meta` element is  
244 added, for which the package vignette on S4 subsetting of `nexml` contains examples.

245 Because NeXML expresses all metadata using the RDF standard, and stores them compliant with  
246 RDFa, they can be extracted as an RDF graph, queried, analyzed, and mashed up with other RDF data,  
247 local or on the web, using a wealth of off-the-shelf tools for working with RDF (see Prud'hommeaux  
248 (2014) or Hartig (2012)). Examples for these possibilities are included in the RNeXML SPARQL vignette,

249 and the package also comes with a demonstration that can be run from R using the following command:  
250 `demo("sparql", "RNeXML")`).

## 251 *Using metadata to extend the NeXML standard*

252 NeXML was designed to prevent the need for future non-interoperable “flavors” of the standard in  
253 response to new research directions. Its solution to this inevitable problem is a highly flexible metadata  
254 system without sacrificing strict validation of syntax and structure.

255 Here we illustrate how RNeXML’s interface to NeXML’s metadata system can be used to record and  
256 share a type of phylogenetic data not taken into account when NeXML was designed, in this case  
257 stochastic character maps (Huelsenbeck *et al.* 2003). Such data assign certain parts (corresponding to  
258 time) of each branch in a time-calibrated phylogeny to a particular “state” (typically of a morphological  
259 characteristic). The current de-facto format for sharing stochastic character maps, created by `simmap`  
260 (Bollback 2006), a widely used tool for creating such maps, is a non-interoperable modification of the  
261 standard Newick tree format.

262 To express these data as a custom metadata extension, we create metadata annotations for the `edge`  
263 elements in the phylogeny topology. The metadata need to describe the character state being assigned,  
264 and the duration (in terms of branch length) that the edge spends in that state. The NeXML standard  
265 already defines `state` elements for discrete character traits, which, as most elements in NeXML, also  
266 have unique identifiers by which they can be referenced.

267 In the following proposal, we term each series of character state assignments for an edge a *reconstruction*,  
268 and for each state assignment comprising a reconstruction, we provide the length, the identifier of the  
269 character state, and an ordinal property to record the chronological order of the assignment. In this  
270 manner, a single edge containing one state change could be annotated as follows. (the identifier values  
271 in the proposal, “cr1”, “s1”, and “s2”, are identifiers local to a – fictitious – NeXML document, based  
272 on the example `simmap` tree illustrated in Revell (2012)):

```
m <- meta("simmap:reconstructions", children = c(  
  meta("simmap:reconstruction", children = c(  
  
    meta("simmap:char", "cr1"),  
    meta("simmap:stateChange", children = c(  

```

```

    meta("simmap:order", 1),
    meta("simmap:length", "0.2030"),
    meta("simmap:state", "s2"))),

    meta("simmap:char", "cr1"),
    meta("simmap:stateChange", children = c(
      meta("simmap:order", 2),
      meta("simmap:length", "0.0022"),
      meta("simmap:state", "s1")))
  ))))

```

273 The example uses a prefix, `simmap`, to group the newly introduced metadata properties in a vocabulary,  
 274 for which the `add_namespace()` method can be used to give a URI as an identifier:

```
nex <- add_namespaces(c(simmap = "https://github.com/ropensci/RNeXML/tree/master/inst/simmap.md")
```

275 Here the URI does not resolve to a fully machine-readable definition of the terms and their semantics,  
 276 but it can nonetheless be used to provide at least a human-readable informal definition of the terms.  
 277 Because the above representation does not modify the NeXML format, it should be readable – even if  
 278 not necessarily interpretable as a stochastic character mapping – by any compliant NeXML parser. In  
 279 addition, it includes a direct and machine-readable link to metadata documentation.

280 Writing out metadata annotations as in the above proposal may seem tedious. However, by virtue of  
 281 RNeXML’s API it is straightforward to define functions that encapsulate these details. As a demonstration,  
 282 RNeXML comes with functions `simmap_to_nexml()` and `nexml_to_simmap()` that map bidirectionally  
 283 between NeXML files with `simmap` metadata according to the above proposal, and the extension of the  
 284 `ape::phylo` class devised by Revell (2012) in the R package `phytools`. The RNeXML `simmap` vignette  
 285 illustrates how this allows seamlessly generating and reading NeXML files containing `simmap` data, and  
 286 using the various functions from `phytools` designed for `simmap` objects (Revell 2012).

## 287 *Publishing NeXML files from R*

288 Data archiving is increasingly required by scientific journals, including in evolutionary biology, ecology,  
 289 and biodiversity (e.g. Rausher et al. (2010)). The effort involved with preparing and submitting

properly annotated data to archives remains a notable barrier to the broad adoption of data archiving and sharing as a normal part of the scholarly publication workflow (Tenopir *et al.* 2011; Stodden 2014). In particular, the majority of phylogenetic trees published in the scholarly record are inaccessible or lost to the research community (Drew *et al.* 2013).

One of RNeXML’s aims is to promote the archival of well-documented phylogenetic data in scientific data repositories, in the form of NeXML files. To this end, the method `nexml_publish()` provides an API directly from within R that allows data archival to become a step programmed into data management scripts. Initially, the method supports the data repository Figshare (<http://figshare.com>):

```
doi <- nexml_publish(birds, visibility = "public", repository="figshare")
```

Figshare also supports “private” visibility, allowing to securely backup data to a repository and to share them with collaborators prior to public release.

## Conclusions and future directions

RNeXML allows R’s ecosystem to read and write data in the NeXML format through an interface that is no more involved than reading or writing data from other phylogenetic data formats. It also carries immediate benefits for its users compared to other formats. For example, comparative analysis R packages and users frequently add their own metadata annotations to the phylogenies they work with, such as annotations of species, stochastic character maps, trait values, model estimates and parameter values. RNeXML affords R the capability to harness machine-readable semantics and an extensible metadata schema to capture, preserve, and share these and other kinds of information, all through an API instead of having to understand in detail the schema underlying the NeXML standard. To assist users in meeting the rising bar for best practices in data sharing in phylogenetic research (Cranston *et al.* 2014), RNeXML captures metadata information from the R environment to the extent possible, and applies reasonable defaults.

The goals for continued development of RNeXML revolve primarily around better interoperability with other existing phylogenetic data representations in R, such those found in the `phylobase` package (NESCENT R Hackathon Team 2014); and better integration of the rich metadata semantics found in ontologies defined in the Web Ontology Language (OWL), including programmatic access to machine reasoning with such metadata.

## 317 Acknowledgements

318 This project was supported in part by the National Evolutionary Synthesis Center (NESCent) (NSF  
319 #EF-0905606), and grants from the National Science Foundation (DBI-1306697) and the Alfred P Sloan  
320 Foundation (Grant 2013-6-22). RNeXML started as a project idea for the Google Summer of Code(TM),  
321 and we thank Kseniia Shumelchik for taking the first steps to implement it. We are grateful to F.  
322 Michonneau for helpful comments on an earlier version of this manuscript.

## 323 References

- 324 Bollback, J. (2006). *BMC Bioinformatics*, **7**, 88. Retrieved from [http://dx.doi.org/10.1186/1471-2105-7-](http://dx.doi.org/10.1186/1471-2105-7-88)  
325 [88](http://dx.doi.org/10.1186/1471-2105-7-88)
- 326 Chamberlain, S.A. & Szöcs, E. (2013). taxize: taxonomic search and retrieval in R. *F1000Research*.  
327 Retrieved from <http://dx.doi.org/10.12688/f1000research.2-191.v2>
- 328 Cranston, K., Harmon, L.J., O’Leary, M.A. & Lisle, C. (2014). Best Practices for Data Shar-  
329 ing in Phylogenetic Research. *PLoS Curr*. Retrieved from [http://dx.doi.org/10.1371/currents.tol.](http://dx.doi.org/10.1371/currents.tol.bf01eff4a6b60ca4825c69293dc59645)  
330 [bf01eff4a6b60ca4825c69293dc59645](http://dx.doi.org/10.1371/currents.tol.bf01eff4a6b60ca4825c69293dc59645)
- 331 Drew, B.T., Gazis, R., Cabezas, P., Swithers, K.S., Deng, J., Rodriguez, R., Katz, L.A., Crandall,  
332 K.A., Hibbett, D.S. & Soltis, D.E. (2013). Lost Branches on the Tree of Life. *PLoS Biol*, **11**, e1001636.  
333 Retrieved from <http://dx.doi.org/10.1371/journal.pbio.1001636>
- 334 Harmon, L., Weir, J., Brock, C., Glor, R. & Challenger, W. (2008). GEIGER: investigating evolutionary  
335 radiations. *Bioinformatics*, **24**, 129–131.
- 336 Hartig, O. (2012). An Introduction to SPARQL and Queries over Linked Data. *Web Engineering* pp.  
337 506–507. Springer Science + Business Media. Retrieved from [http://dx.doi.org/10.1007/978-3-642-](http://dx.doi.org/10.1007/978-3-642-31753-8_56)  
338 [31753-8\\_56](http://dx.doi.org/10.1007/978-3-642-31753-8_56)
- 339 Huelsenbeck, J.P., Nielsen, R. & Bollback, J.P. (2003). Stochastic Mapping of Morphological Characters.  
340 *Systematic Biology*, **52**, 131–158. Retrieved from <http://dx.doi.org/10.1080/10635150390192780>
- 341 Lang, D.T. (2013). *XML: Tools for parsing and generating XML within R and S-Plus*. Retrieved from  
342 <http://CRAN.R-project.org/package=XML>
- 343 Maddison, D., Swofford, D. & Maddison, W. (1997). NEXUS: an extensible file format for systematic  
344 information. *Syst. Biol.*, **46**, 590–621. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/11975335>

345 NESCENT R Hackathon Team. (2014). *phylobase: Base package for phylogenetic structures and*  
346 *comparative data*. Retrieved from <http://CRAN.R-project.org/package=phylobase>

347 O'Meara, B. (2014). CRAN Task View: Phylogenetics, Especially Comparative Methods. Retrieved  
348 from <http://cran.r-project.org/web/views/Phylogenetics.html>

349 Paradis, E., Claude, J. & Strimmer, K. (2004). APE: analyses of phylogenetics and evolution in R  
350 language. *Bioinformatics*, **20**, 289–290.

351 Parr, C.S., Guralnick, R., Cellinese, N. & Page, R.D.M. (2011). Evolutionary informatics: unifying  
352 knowledge about the diversity of life. *Trends in ecology & evolution*, **27**, 94–103. Retrieved from  
353 <http://www.ncbi.nlm.nih.gov/pubmed/22154516>

354 Piel, W.H., Donoghue, M.J. & Sanderson, M.J. (2002). TreeBASE: a database of phylogenetic  
355 information. *The Interoperable 'Catalog of Life'* (eds J. Shimura, K.L. Wilson & D. Gordon), pp.  
356 41–47. Research Report. National Institute for Environmental Studies, Tsukuba, Japan. Retrieved  
357 from [http://donoghuelab.yale.edu/sites/default/files/124\\_piel\\_shimura02.pdf](http://donoghuelab.yale.edu/sites/default/files/124_piel_shimura02.pdf)

358 Prud'hommeaux, E. (2014). SPARQL Query Language for RDF. *W3C*. Retrieved from <http://www.w3.org/TR/rdf-sparql-query/>

360 R Core Team. (2014). *R: A Language and Environment for Statistical Computing*. R Foundation for  
361 Statistical Computing, Vienna, Austria. Retrieved from <http://www.R-project.org/>

362 Rausher, M.D., McPeck, M.A., Moore, A.J., Rieseberg, L. & Whitlock, M.C. (2010). Data Archiving.  
363 *Evolution*, **64**, 603–604. Retrieved from <http://dx.doi.org/10.1111/j.1558-5646.2009.00940.x>

364 Revell, L.J. (2012). phytools: An R package for phylogenetic comparative biology (and other things).  
365 *Methods in Ecology and Evolution*, **3**, 217–223.

366 Stodden, V. (2014). The Scientific Method in Practice: Reproducibility in the Computational Sciences.  
367 *SSRN Journal*. Retrieved from <http://dx.doi.org/10.2139/ssrn.1550193>

368 Stoltzfus, A., O'Meara, B., Whitacre, J., Mounce, R., Gillespie, E.L., Kumar, S., Rosauer, D.F. & Vos,  
369 R.A. (2012). Sharing and re-use of phylogenetic trees (and associated data) to facilitate synthesis. *BMC*  
370 *Research Notes*, **5**, 574. Retrieved from <http://dx.doi.org/10.1186/1756-0500-5-574>

371 Tenopir, C., Allard, S., Douglass, K., Aydinoglu, A.U., Wu, L., Read, E., Manoff, M. & Frame, M.  
372 (2011). Data Sharing by Scientists: Practices and Perceptions (C. Neylon, Ed.). *PLoS ONE*, **6**, e21101.  
373 Retrieved from <http://dx.doi.org/10.1371/journal.pone.0021101>

374 Vos, R.A., Balhoff, J.P., Caravas, J.A., Holder, M.T., Lapp, H., Maddison, W.P., Midford, P.E.,  
375 Priyam, A., Sukumaran, J., Xia, X. & Stoltzfus, A. (2012). NeXML: Rich, Extensible, and Verifiable  
376 Representation of Comparative Data and Metadata. *Systematic Biology*, **61**, 675–689. Retrieved from  
377 <http://dx.doi.org/10.1093/sysbio/sys025>