1 **Title:** helminthR: An R interface to the London Natural History Museum's Host-Parasite
2 Database
3 **Running title:** Global helminth data access
4 **Author:** Tad Dallas
5 **Affiliation**: University of Georgia, Odum School of Ecology, Athens Ga 30602
6 **Email:** tdallas@uga.edu

7

## 8 **Abstract**

9 The understanding of the diversity and distribution of helminth parasites is currently constrained

10 by the limited number of host-parasite interaction databases, and the difficulty in accessing

11 existing data. The London Natural History Museum's Host-Parasite Database represents one

12 such underutilized database, containing over a quarter million helminth parasite occurrence

13 records, accessible through a web interface. To enable users to programmatically search and

14 manipulate data from this database, I developed an R package called `helminthR`. Here, I

15 introduce the core functions of the package, and detail how `helminthR` can be used to obtain

16 host-parasite interaction records, citations for interactions, and host taxonomic data.

17


18
19
20
21
22
23
24

## Introduction

Helminth parasites are one of the most common infectious agents to humans (Stoll 1947, Hotez et al. 2008, De Silva et al. 2003), wild animals (Poulin and Valtonen 2002, Jolles et al. 2008), and livestock (Over et al. 1992, Morgan et al. 2013). Limitations in data availability have hampered our understanding of the spatial distribution of helminth parasites, and associations between helminth parasites and both human and wildlife hosts. Further, there is a need for basic scientific research into the community ecology and macroecology of host-helminth associations (Rohde 2002). Such efforts could provide tests of principles from community ecology, and macroecological patterns in parasites.

To address these research concerns, data on host-helminth associations across broad spatial scales are needed. Efforts to document known host-parasite associations in large databases are fairly recent, and represent valuable resources for researchers (see Strona et al. 2013, Gibson et al. 2005, Nunn and Altizer 2005). However, a portion of these databases are not openly accessible, requiring users to contact database administrators or to copy data from web interfaces. These methods of accessing databases may lead to transcription errors, duplicated efforts among labs, and create static copies of the data that are difficult to update if and when new data are added. Allowing host-parasite databases to be open and easy to access may promote open and reproducible science, and would potentially promote the discovery of "general laws" in parasite ecology (Poulin 2007).

47　To this end, I have developed an R package capable of extracting information from a large global

48　database of host-helminth parasite occurrence records maintained by the London Natural History

49　Museum (NHM; Gibson et al. 2005). This curated database includes more than 250,000 host-

50　helminth records from over 28,000 published peer-reviewed articles. However, the web interface

51　of the database makes data analysis difficult, which subsequently limits the use of this data

52　resource by researchers (but see Strona and Fattorini (2014) and Wells et al. (2015)). The goal of

53　the `helminthR` package is to make all the data contained in the London Natural History

54　Museum's database accessible from R, a commonly used open source statistical programming

55　environment (R core team 2015).

## 56　Core package functionality

57

58　Here, I explore the core functions of the `helminthR` package, and then demonstrate the utility of

59　`helminthR` for creating host-parasite interaction networks. `helminthR` relies on several packages

60　that interface with html and xml, including `rvest` (Wickham 2015) and `xml2` (Wickham 2015b).

61　Currently, `helminthR` is available on Github, and is hosted by the rOpenSci collective, a group

62　of scientists and developers committed to creating packages to promote open science, including

63　the creation of packages to access online data sources. The package can be easily downloaded

64　using the `devtools` package, using the following R code.

65

```
66    devtools::install_github('ropensci/helminthR')
67    library('helminthR')
```
68

69

70  Downloading and using this package does not require the user to have a Github account, unless

71  they would like to actively contribute to package functionality, or file an issue.

72

73  ***Querying the database***

74  Host-parasite records in the NHM database contain information on host and parasite species, one

75  or more citations for the host-parasite association, and the location of the interaction

76  georeferenced to the country, state (for the United States), or water body (e.g. Lake Erie) level.

77  Queries can be made to find all interactions of a known host species (`findHost`), all interactions

78  of a known parasite species (`findParasite`), or all interactions at a specific geographic location

79  (`findLocation`). Links to citations for a given helminth record can be obtained from any of the

80  functions listed above by setting the `citation` argument to `TRUE`.

81

82  When querying the database for known hosts or helminths, the user can input genus and/or

83  species name in order to query different taxonomic levels of host or parasite. Further,

84  `findParasite` can find host-helminth records given a parasite group (Cestodes,

85  Acanthocephalans, Monogeneans, Nematodes, Trematodes, or Turbellarian) or subgroup. The

86  following example code would find all interactions of nematodes in the genus *Strongyloides*.

87
88  
```
  StrongHosts <- findParasite(genus='Strongyloides', validateHosts=FALSE)
```
89

90
91
92  The resulting structure of `strongHosts` is a host-parasite matrix in the form of a three (or four)

93  column `data.frame` containing host and parasite names, parasite full name, and citation (if the

94  `citation` argument is set as `TRUE`). The argument `validateHosts` provides taxonomic

4

95  information on hosts from the Catalogue of Life (Roskov et al. 2015). While slightly slow, this

96  removes questionable hosts, and validates species names (when `validateHosts=TRUE`),

97  returning a `list` object containing the `data.frame` described above, and the taxonomic

98  information for all hosts. This structure is maintained when querying using any of the "find*"

99  functions, including `findHost, findParasite,` and `findLocation`. The following code

100 demonstrates the `findHost` function in order to find helminth occurrence records in wild

101 individuals of *Gorilla gorilla* (using the `hostState` argument). The user can also query captive

102 hosts, domesticated hosts, or hosts used in commercial applications.

103

```
104     gorillaParasites <- findHost(genus = 'Gorilla', species = 'gorilla',
105                                    hostState = 1 )
106
107
108
```

109 The final core function in the `helminthR` package queries all host-parasite interactions for a

110 given geographic location. A list of locations capable of being queried is provided by the

111 `listLocations` function, and a cached copy of these data is provided as a data object (using the

112 command `data(locations)`). Georeferencing of these data is performed using the `geocode`

113 function in the `ggmap` package (Kahle & Wickham 2013). The user is responsible for ensuring

114 the accuracy of the provided latitude and longitude coordinates. Further care should be taken

115 when searching by location, as some locations may be nested within others (e.g. "South

116 America" is a valid location query, but many countries in South America are also valid queries).

117 Below, I demonstrate the functionality by finding all host-parasite associations recorded in

118 France where the host was "in the wild" (i.e., `hostState = 1`), removing occurrence records

119 where the host or parasite has parantheses (e.g. "(freshwater_fish)") or is identified to be at the

120 genus level (e.g. "Sanguinicola spp.") by setting the argument `speciesOnly` to be `TRUE`. The

121 result is a host-parasite association list containing information on host-helminth associations,

122 including links to the original citations. It is important to note that not all interactions will be

123 unique, so the user must use the `unique` function on the `Host` and `Parasite` columns of the

124 output `data.frame`.

125

```
126     # Find all host-helminth associations occuring in France
127     FrenchHostPars <- findLocation(location = 'France', speciesOnly = TRUE,
128                                    citation = TRUE)
129
130     # Find unique host-parasite associations
131     FrenchHostParsUnique <- unique(FrenchHostPars[,1:2])
132
```

133 *Visualizing host-parasite networks*

134 The above code demonstrates the functionality of the `helminthR` package for querying host-

135 parasite interactions by host and parasite genus and/or species, and also for locating all host-

136 parasite interactions in a given country or locality. Using the `findLocation` function, I queried

137 the database for all host-parasite interactions occurring within Lake Erie, one of the US Great

138 Lakes, and visualized the resulting host-parasite interaction network (Figure 1) using the `igraph`

139 `R` package (Csardi & Nepusz 2006). Detailed code to create this type of visualization is provided

140 in the supplement.

141 **Data limitations**

142

143 The data contained in the London Natural History Museum's Host-Parasite Database represent a

144 valuable resource, but are not without limitation. First, the data are from studies published

145 anytime after 1922, and the data owners themselves accept no responsibility for data accuracy.

146 Second, the data are only georeferenced to the country level in most cases, which limits their

147 application. However, citations are given for each host-parasite association, and an attempt has

148 been made to obtain latitude and longitude values for the centroids of countries (using the

149 command `data(locations)`). While this may be time consuming, the examination of original

150 references would help assure data quality, and provide more fine georeferencing. Nevertheless,

151 the data can still be used to address many macroecological patterns in their current form.  For

152 example, data on aquatic and marine parasites are georeferenced to coastal areas (e.g. "Coast of

153 New Guinea") or larger bodies of water (e.g. "Aral sea"), providing a way to apply

154 macroecological theory to largely unexplored questions related to the diversity and distribution

155 of marine parasites (Rohde 2002, Rhode 2010).


156 **Conclusions**

157

158 In this paper I have shown how the R package `helminthR` permits the programmatic access of

159 the Natural History Museum Host-Parasite Database, making it easy to generate host-parasite

160 networks at different geographical scales spanning from local to global. This database represents

161 one of the most complete aquatic host-parasite databases (but see Strona et al. 2013), providing

162 data on parasite occurrences for both terrestrial and aquatic hosts. With any luck, `helminthR` will

163 promote the application of concepts from community ecology and macroecology to parasite

164 communities at a broader spatial scale. This project is hosted on Github, and uses TravisCI for

165 continuous integration of the package on different R versions. Issues or improvements can be

166 suggested at this link (https://github.com/ropensci/helminthR/issues).

167

168  To cite `helminthR` or acknowledge its use, cite the original data source (Gibson et al. 2005), and

169  this Software note as follows, substituting the version of the application that you used for `ver.

170  xxx`:

171  T. Dallas 2015. `helminthR`: An R interface to the London Natural History Museum's Host-

172  Parasite Database - Ecography (ver. xxx).

## Acknowledgments
173

174

180

## References
181

182

183  Crompton, D. W. T. and Nesheim, M. 2002. Nutritional impact of intestinal helminthiasis during

184  the human life cycle.  -Annual Review of Nutrition 22: 35-59.

185

186  Csardi G, Nepusz T. 2006.  The igraph software package for complex network research.

187  -InterJournal, Complex Systems 1695. http://igraph.org

188

189   De Silva, N. R. et al. 2003. Soil-transmitted helminth infections: updating the global picture.

190   -Trends in Parasitology 19: 547-551.

191

192   Gibson, D. et al. 2005. Host-parasite database of the natural history museum, London.

193   http://www.nhm.ac.uk/research-curation/scientific-resources/taxonomy-systematics/host-

194   parasites/database/index.jsp .

195

196   Hotez, P. J. et al. 2008. Helminth infections: the great neglected tropical diseases.  -Journal of

197   Clinical Investigation 118: 1311-1321.

198

199   Jolles, A. E. et al. 2008. Interactions between macroparasites and microparasites drive infection

200   patterns in free-ranging african buffalo.  -Ecology 89: 2239-2250.

201

202   Kahle, D. and H. Wickham. 2013. ggmap: Spatial Visualization with ggplot2. -The R Journal,

203   5(1): 144-161. http://journal.r-project.org/archive/2013-1/kahle-wickham.pdf

204

205   Morgan, E. R. et al. 2013. Global change and helminth infections in grazing ruminants in

206   Europe: impacts, trends and sustainable solutions. -Agriculture 3: 484-502.

207

208   Nunn, C. L. and Altizer, S. M. 2005. The global mammal parasite database: an online resource

209   for infectious disease records in wild primates.  -Evolutionary Anthropology: Issues, News, and

210   Reviews 14: 1-2.

211

212  Over, H. J. et al. 1992. Distribution and impact of helminth diseases of livestock in developing

213  countries. 96. -Food & Agriculture Organization.

214

215  Poulin, R. and Valtonen, E. T. 2002. The predictability of helminth community structure in space:

216  a comparison of fish populations from adjacent lakes. -International Journal for Parasitology 32:

217  1235-1243.

218

219  Poulin, R. 2007. Are there general laws in parasite ecology?.-Parasitology 134(06): 763-776.

220

221  R Core Team 2015. R: A Language and Environment for Statistical Computing. R Foundation for

222  Statistical Computing, Vienna, Austria.

223

224  Rohde, K. 2002. Ecology and biogeography of marine parasites. - Advances in Marine Biology

225  43: 1-83.

226

227  Rohde, K. 2010. Marine parasite diversity and environmental gradients. In: Morand, S, and BR

228  Krasnov, (eds.) The biogeography of host-parasite interactions. Oxford, UK: Oxford University

229  Press. pp. 73-88.

230

231     Roskov Y. et. al. 2015. Species 2000 & ITIS Catalogue of Life, 2015 Annual Checklist. Digital

232     resource at www.catalogueoflife.org/annual-checklist/2015. Species 2000: Naturalis, Leiden, the

233     Netherlands. ISSN 2405-884X.

234

235     Stephenson, L. S. et al. 2000. Malnutrition and parasitic helminth infections. -Parasitology 121:

236     S23-S38.

237

238     Stoll, N. R. 1947. This wormy world. -The Journal of Parasitology 33: 1.

239

240     Strona, G. and Fattorini, S. 2014. Parasitic worms: how many really?  -International Journal for

241     Parasitology 44: 269-272.

242

243     Strona, G. et al. 2013. Host range, host ecology, and distribution of more than 11,800 fish

244     parasite species: Ecological archives e094-045. -Ecology 94: 544-544.

245

246     Wells, K. et al. 2015. The importance of parasite geography and spillover effects for global

247     patterns of host -parasite associations in two invasive species. -Diversity and Distributions

248     14 21: 477-486.

249

250     Wickham, H. 2015. rvest: Easily Harvest (Scrape) Web Pages. R  package version 0.3.1.

251     http://CRAN.R-project.org/package=rvest

252

253    Wickham, H. 2015b. xml2: Parse XML. R package version 0.1.2. http://CRAN.R-

254    project.org/package=xml2

255

256    **Figure captions**

257

258    FIG. 1. The host-parasite association network for Lake Erie, one of the Great Lakes located in the

259    Northern United States. Grey lines between boxes represent interactions between hosts (larger

260    blue dots) and helminth parasites (smaller black dots).

261