

mangal – making ecological network analysis simple

Timothée Poisot et al. (see complete author list below)

Author list: Timothée Poisot (1,2), Benjamin Baiser (3), Jennifer A. Dunne (4,5), Sonia Kéfi (6), François Massol (7,8), Nicolas Mouquet (6), Tamara N. Romanuk (9), Daniel B. Stouffer (1), Spencer A. Wood (10,11), Dominique Gravel (12,2)

1. University of Canterbury, School of Biological Sciences, Christchurch, New Zealand
2. Québec Centre for Biodiversity Sciences, Montréal (QC), Canada
3. Department of Wildlife Ecology and Conservation, University of Florida, Gainesville
4. Sante Fe Institute, 1399 Hyde Park Road, Santa Fe NM 87501
5. Pacific Ecoinformatics and Computational Ecology Lab, 1604 McGee Ave., Berkeley, CA 94703
6. Institut des Sciences de l'Évolution, UMR CRNS 5554, Université Montpellier 2, 3405 Montpellier, France
7. Laboratoire Génétique et Evolution des Populations Végétales, CNRS UMR 8198, Université Lille 1, Bâtiment SN2, F-59655 Villeneuve d'Ascq cedex, France
8. UMR 5175 CEFÉ – Centre d'Ecologie Fonctionnelle et Evolutive (CNRS), 1919 Route de Mende, F-34293 Montpellier cedex 05, France
9. Department of Biology, Dalhousie University
10. Natural Capital Project, School of Environmental and Forest Sciences, University of Washington, Seattle, WA 98195, USA
11. Department of Biological Sciences, Idaho State University, Pocatello, ID 83209, USA
12. Université du Québec à Rimouski, Département de Biologie, 300 Allées des Ursulines, Rimouski (QC) G5L 3A1, Canada

Author for correspondence: tim@poisotlab.io.

1 **Short title:** Automated retrieval of ecological interaction data

2 **Keywords:** R, API, database, open data, ecological networks, species interactions

3 The study of ecological networks is severely limited by (i) the difficulty to access data, (ii) the lack of a
4 standardized way to link meta-data with interactions, and (iii) the disparity of formats in which ecological
5 networks themselves are stored and represented. To overcome these limitations, we have designed a data
6 specification for ecological networks. We implemented a database respecting this standard, and released
7 a R package (`rmangal`) allowing users to programmatically access, curate, and deposit data on ecological
8 interactions. In this article, we show how these tools, in conjunction with other frameworks for the program-
9 matic manipulation of open ecological data, streamlines the analysis process and improves replicability and
10 reproducibility of ecological network studies.

1 Introduction

2 Ecological networks are efficient representations of the complexity of natural communities, and help discover mechanisms
3 contributing to their persistence, stability, resilience, and functioning. Most of the early studies of ecological networks
4 were focused on understanding how the structure of interactions within one location affected the ecological properties
5 of this local community. They revealed the contribution of average network properties, such as the buffering impact
6 of modularity on species loss (Pimm et al. 1991,??), the increase in robustness to extinctions along with increases in
7 connectance (Dunne et al. 2002), and the fact that organization of interactions maximizes biodiversity (Bastolla et al.
8 2009). New studies introduced the idea that networks can vary from one locality to another. They can be meaningfully
9 compared, either to understand the importance of environmental gradients on the presence of ecological interactions
10 (Tylianakis et al. 2007), or to understand the mechanisms behind variation itself (Poisot et al. 2012, 2014). Yet, meta-
11 analyses of numerous ecological networks are still extremely rare, and most of the studies comparing several networks do
12 so within the limit of particular systems (Schleuning et al. 2011, Dalsgaard et al. 2013, Poisot et al. 2013, Chamberlain
13 et al. 2014, Olito and Fox 2014). The severe shortage of publicly shared data in the field also restricts the scope of
14 large-scale analyses.

15 It is possible to predict the structure of ecological networks, either using latent variables (Rohr et al. 2010, Eklöf et al.
16 2013) or actual trait values (Gravel et al. 2013). The calibration of these approaches require accessible data, not only
17 about the interactions, but about the traits of the species involved. Comparing the efficiency of different methods is also
18 facilitated if there is a homogeneous way of representing ecological interactions, and the associated metadata. In this
19 paper, we (i) establish the need of a data specification serving as a common language among network ecologists, (ii)
20 describe this data specification, and (iii) describe `rmangal`, a R package and companion database relying on this data
21 specification. The `rmangal` package allows to easily deposit and retrieve data about ecological interactions and networks
22 in a publicly accessible database. We provide use cases showing how this new approach makes complex analyzes simpler,
23 and allows for the integration of new tools to manipulate biodiversity resources.

24 Networks need a data specification

25 Ecological networks are (often) stored as an *adjacency matrix* (or as the quantitative link matrix), that is a series of 0s
26 and 1s indicating, respectively, the absence and presence of an interaction. This format is extremely convenient for *use*
27 (as most network analysis packages, *e.g.* `bipartite`, `betalink`, `foodweb`, require data to be presented this way), but
28 is extremely inefficient at *storing* meta-data. In most cases, an adjacency matrix provides information about the identity
29 of species (in the cases where rows and columns headers are present) and the presence or absence of interactions. If
30 other data about the environment (*e.g.* where the network was sampled) or the species (*e.g.* the population size, trait
31 distribution, or other observations) are available, they are often either given in other files or as accompanying text. In both
32 cases, making a programmatic link between interaction data and relevant meta-data is difficult and, more importantly,
33 error-prone.

34 By contrast, a data specification (*i.e.* a set of precise instructions detailing how each object should be represented) provides
35 a common language for network ecologists to interact, and ensures that, regardless of their source, data can be used in
36 a shared workflow. Most importantly, a data specification describes how data are *exchanged*. Each group retains the
37 ability to store the data in the format that is most convenient for in-house use, and only needs to provide export options
38 (*e.g.* through an API, *i.e.* a programmatic interface running on a web server, returning data in response to queries in
39 a pre-determined language) respecting the data specification. This approach ensures that *all* data can be used in meta-

analyses, and increases the impact of data (Piwowar and Vision 2013). Data archival also offers additional advantages for ecology. The aggregation of local observations can reveal large-scale phenomena (Reichman et al. 2011), which would be unattainable in the absence of a collaborative effort. Data archival in databases also prevents data rot and data loss (Vines et al. 2014), thus ensuring that data on interaction networks – which are typically hard and costly to produce – continue to be available and usable.

6 Elements of the data specification

The data specification introduced here (Fig. 1) is built around the idea that (ecological) networks are collections of relationships between ecological objects, and each element has particular meta-data associated with it. In this section, we detail the way networks are represented in the `mangal` specification. An interactive webpage with the elements of the data specification can be found online at <http://mangal.io/doc/spec/>. The data specification is available either at the API root (e.g. <http://mangal.io/api/v1/?format=json>), or can be viewed using the `whatIs` function from the `rmangal` package. Rather than giving an exhaustive list of the data specification (which is available online at the aforementioned URL), this section serves as an overview of each element, and how they interact.

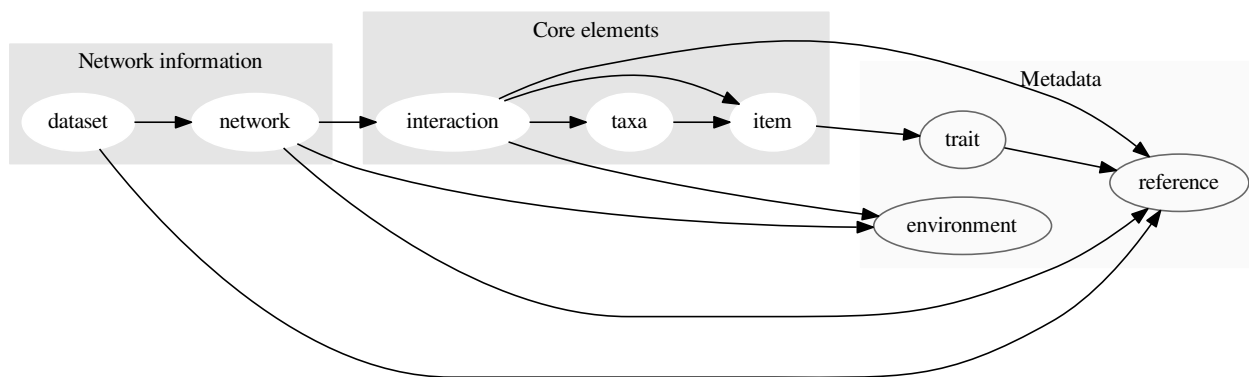


Fig. 1: An overview of the data specification, and the hierarchy between objects. Every box corresponds to a level of the data specification. Grey boxes are nodes, blue boxes are interactions and networks, and green boxes are metadata. The **bold** boxes (dataset, network, interaction, taxa) are the minimal elements needed to represent a network.

We propose JSON, a user-friendly format equivalent to XML, as an efficient way to standardise data representation for two main reasons. First, it has emerged as a *de facto* standard for web platform serving data, and accepting data from users. Second, it allows strict *validation* of the data: a JSON file can be matched against a scheme, and one can verify that it is correctly formatted (this includes the possibility that not all fields are filled, as will depend on available data). Finally, JSON objects are easily and cheaply (memory-wise) parsed in the most commonly-used programming languages, notably R (equivalent to `list`) and python (equivalent to `dict`). For most users, the format in which data are transmitted is unimportant, as the interaction happens within R – as such, knowing how JSON objects are organized is only useful for those who want to interact with the API directly. As such, the `rmangal` package takes care of converting the data into the correct JSON format to upload them in the database.

Functions in the `rmangal` package are names after elements of the data specification, in the following way: `verb +`

1 `Element`. `verb` can be one of `list`, `get`, or `patch`; for example, the function to get a particular network is `getNetwork`.
2 The function to modify (patch) a taxon is `patchTaxa`. All of these functions return a *list* object, which means that chaining
3 them together using, *e.g.* the `plyr` package, is time-efficient. There are examples of this in the use-cases.

4 **Node information**

5 **Taxa**

6 Taxa are a taxonomic entity of any level, identified by their name, vernacular name, and their identifiers in a variety of
7 taxonomic services. Associating the identifiers of each taxa allows using the new generation of open data tools, such
8 as `taxize` (Chamberlain and Szöcs 2013), in addition to protecting the database against taxonomic revisions. The data
9 specification currently has fields for `ncbi` (National Center for Biotechnology Information), `gbif` (Global Biodiversity
10 Information Facility), `tsn` (Taxonomic Serial Number, used by the Integrated Taxonomic Information System), `eol`
11 (Encyclopedia of Life) and `bold` (Barcode of Life) identifiers. We also provide the taxonomic status, *i.e.* whether the
12 taxon is a true taxonomic entity, a “trophic species”, or a morphospecies. Taxonomic identifiers can either be added by
13 the contributors, or will be automatically retrieved during the automated curation routine.

14 **Item**

15 An `item` is any measured instance of a taxon. Items have a `level` argument, which can be either `individual` or
16 `population`; this allows representing both individual-level networks (*i.e.* there are as many `items` of a given taxa as
17 there were individuals of this taxon sampled), and population-level networks. When `item` represents a population, it
18 is possible to give a measure of the size of this population. The notion of `item` is particularly useful for time-replicated
19 designs: each observation of a population at a time-point is an `item` with associated `trait` values, and possibly population
20 size.

21 **Network information**

22 All objects described in this sub-section can have a spatial position, information on the date of sampling, and references
23 to both papers and datasets.

24 **Interaction**

25 An `interaction` links two `taxa` objects (but can also link pairs of `items`). The most important attributes of `interactions`
26 are the type of interaction (of which we provide a list of possible values, see *Supp. Mat. 1*), and its `ob_type`, *i.e.* how
27 it was observed. This field helps differentiate direct observations, text mining, and inference. Note that the `obs_type`
28 field can also take `confirmed absence` as a value; this is useful for, *e.g.*, “cafeteria” experiments in which there is high
29 confidence that the interaction did not happen.

30 **Network**

31 A `network` is a series of `interaction` objects, along with (i) information on its spatial position (provided at the latitude
32 and longitude), (ii) the date of sampling, and (iii) references to measures of environmental conditions.

1 Dataset

2 A dataset is a collection of one or several network(s). Datasets also have a field for data and papers, both of which
3 are references to bibliographic or web resources that describe, respectively, the source of the data and the papers in which
4 these data have been studied. Datasets or networks are the preferred entry point into the resources, although in some cases
5 it can be meaningful to get a list of interactions only.

6 Meta-data

7 Trait value

8 Objects of type `item` can have associated `trait` values. These consist in the description of the trait being measured,
9 the value, and the units in which the measure was taken. As traits may have been measured at a different time and/or
10 location that the interaction was, they have fields for time, latitude and longitude, and references to original publication
11 and original datasets.

12 Environmental condition

13 Environmental conditions are associated to datasets, networks, and interactions objects, to allow for both macro and micro
14 environmental conditions. These are defined by the environmental property measured, its value, and the units. As traits,
15 they have fields for time, latitude and longitude, and references to original publication and original datasets.

16 References

17 References are associated to datasets. They accommodate the DOI, JSON or PubMed identifiers, or a URL. When
18 possible, the DOI is preferred as it offers more potential to interact with other online tools, such as the *CrossRef* API.

19 Use cases

20 In this section, we present use cases using the `rmangal` package for R, to interact with a database implementing this
21 data specification, and serving data through an API (<http://mangal.io/api/v1/>). It is possible for users to deposit
22 data into this database through the R package. Note that data are made available under a *CC-0 Waiver* (???). Detailed
23 information about how to upload data are given in the vignettes and manual of the `rmangal` package. In addition, the
24 `rmangal` package comes with vignettes explaining how users can upload their data into the database through R.

25 The data we use for this example come from Ricciardi et al. (2010). These data were previously available on the *Interac-*
26 *tionWeb DataBase* as a single xls file. We uploaded them in the mangal database at <http://mangal.io/api/v1/dataset/2>.

27 The `rmangal` package can be installed this way:

```
# Prepare the environment
library(devtools)
# This line is needed on some linux distributions
if(getOption('unzip')=='' ) options('unzip' = 'unzip')
```

```
# This installs the rmangal package
install_github('mangal-wg/rmangal')
library(rmangal)
```

- 1 Once `rmangal` is installed and loaded, users can establish a connection to the database this way:

```
mangal_url <- 'http://mangal.io/'
api <- mangalapi(mangal_url)
```

2 Create taxa and add an interaction

- 3 In the first use-case, we will create an interaction between two taxa. We ask of readers *not* to execute this code as it is,
- 4 but rather to use it as a template for their own analyses. A complete, step-by-step guide to upload is given in the vignettes
- 5 of the `rmangal` package. Uploading anything requires a username and API key, which can be obtained at the following
- 6 URL: <http://mangal.io/dashboard/login>. Your API key be generated automatically after registration. You can use
- 7 it to connect to the database securely:

```
api_secure <- mangalapi("http://mangal.io", usr="MyUserName", key="AbcDefIjKL1234")
```

- 8 The first step is to create two taxa objects, with the species that we observed interacting:

```
seal <- list(
  name = "Hydrurga leptonych",
  vernacular = "Leopard seal",
  eol = 328637
)
cod <- list(
  name = "Gadus morhua",
  vernacular = "Atlantic cod"
)
```

- 9 Now, we will send these two objects in the remote database:

```
seal <- addTaxa(api_secure, seal)
cod <- addTaxa(api_secure, cod)
```

- 10 Note that it is suggested to overwrite the local copy of the object, because the database will *always* send back the remote
- 11 copy. This makes the syntax of further addition considerably easier, as we show below. Once the two objects are created,
- 12 we can create an interaction between them:

```
seal_eats_cod <- list(
  taxa_from = seal,
  taxa_to = cod,
```

```

    int_type = "predation",
    obs_type = "observed"
)

```

- 1 Then using the same approach, we can send this information in the remote database:

```

seal_eats_cod <- addInteraction(api_secure, seal_eats_cod)

```

- 2 To create networks, datasets, etc, one needs follow the same procedure, as is explained in the online guide for data
- 3 contributors, available at <http://mangal.io/doc/upload/>.

4 Link-species relationships

- 5 In the first example, we visualize the relationship between the number of species and the number of interactions, which
- 6 Martinez (1992) proposed to be linear (in food webs).

```

library(plyr)
library(igraph)

# Retrieve the dataset of interest
dataset <- getDataset(api, 2)

# Get each network in the dataset as a graph object
graphs <- alply(dataset$networks, 1, function(x) toIgraph(api, x))

# Make a data.frame with the number of links and species
ls <- ldply(graphs, function(x) c(S = length(V(x)), L = length(E(x))))
ls$X1 <- aaply(as.numeric(as.vector(ls$X1)), 1,
              function(x) getNetwork(api, x)$name)

7 ## Error in eval(expr, envir, enclos): client error: (404) Not Found

colnames(ls)[1] <- 'Network'

# Now plot this dataset
source("suppmat/usecase_ls.r")

```

- 8 Getting the data to produce this figure requires less than 10 lines of code. The only information needed is the identifier of
- 9 the network or dataset, which we suggest should be reported in publications as: “These data were deposited in the mangal
- 10 format at <URL>/api/v1/dataset/<ID>” (where <URL> and <ID> are replaced by the corresponding values), preferably
- 11 in the methods, possibly in the acknowledgements. To encourage data sharing and its recognition, we encourage users of
- 12 the database to always cite the original datasets or publications.

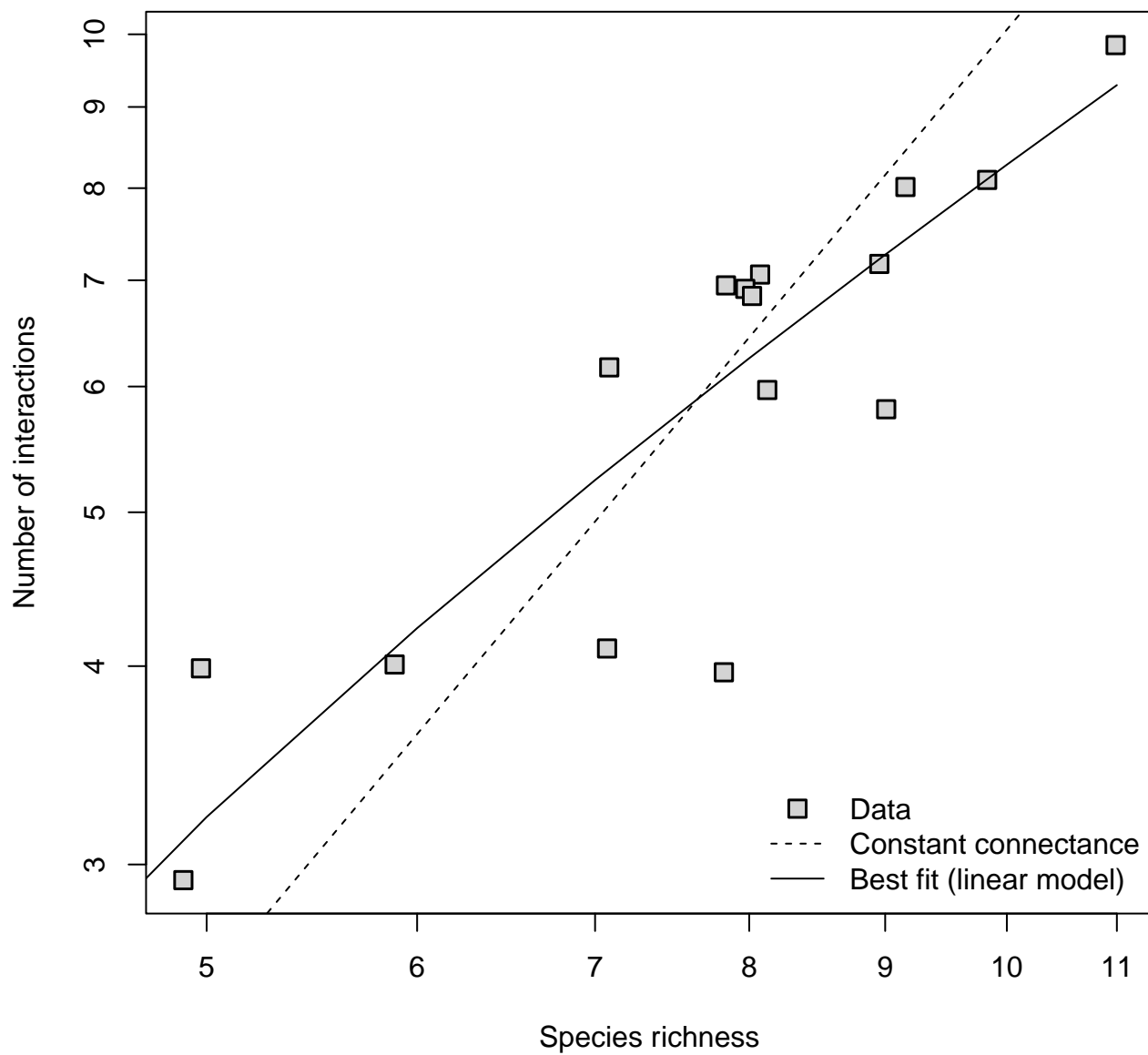


Fig. 2: Relationship between the number of species and number of interactions in the anemonefish-fish dataset. Constant connectance refers to the hypothesis that there is a quadratic relationship between these two quantities.

1 Network beta-diversity

2 In the second example, we use the framework of network β -diversity (Poisot et al. 2012) to measure the extent to which
3 networks that are far apart in space have different interactions. Each network in the dataset has a latitude and longitude,
4 meaning that it is possible to measure the geographic distance between two networks. For each pair of networks, we
5 measure the geographic distance (in km), the species dissimilarity (β_S), the network dissimilarity when all species are
6 present (β_{WN}), and finally, the network dissimilarity when only shared species are considered (β_{OS}).

```
# We need the betalink package to measure network beta-diversity
install_github('tpoisot/betalink')
library(betalink)

library(plyr)
library(igraph)
library(sp)

# We first retrieve all information about the networks
Networks <- alply(dataset$networks, 1, function(x) getNetwork(api, x))

# Extract the lat/lon data
LatLon <- ldply(Networks, function(x) c(name = x$name, lat = x$latitude, lon = x$longitude))
rownames(LatLon) <- LatLon$name
LatLon$lat <- as.numeric(LatLon$lat)
LatLon$lon <- as.numeric(LatLon$lon)
LatLon <- LatLon[,c('lat', 'lon')]

# Then we measure the distances between all pairs of sites
GeoDist <- spDists(as.matrix(LatLon, latlon=TRUE))
colnames(GeoDist) <- rownames(GeoDist) <- rownames(LatLon)
GeoDist <- as.dist(GeoDist)

# Now, we measure the beta-diversity of the networks
names(graphs) <- aapply(names(graphs), 1, function(x) Networks[[x]]$name)
# Finally, we measure the beta-diversity
BetaDiv <- network_betadiversity(graphs)

# We add the geographic distance
BetaDiv$GEO <- GeoDist

# Plotting
source("suppmat/usecase_beta.r")
```

7 As shown in Fig. 3, while species dissimilarity and overall network dissimilarity increase when two networks are far
8 apart, this is not the case for the way common species interact. This suggests that in this system, network dissimilarity

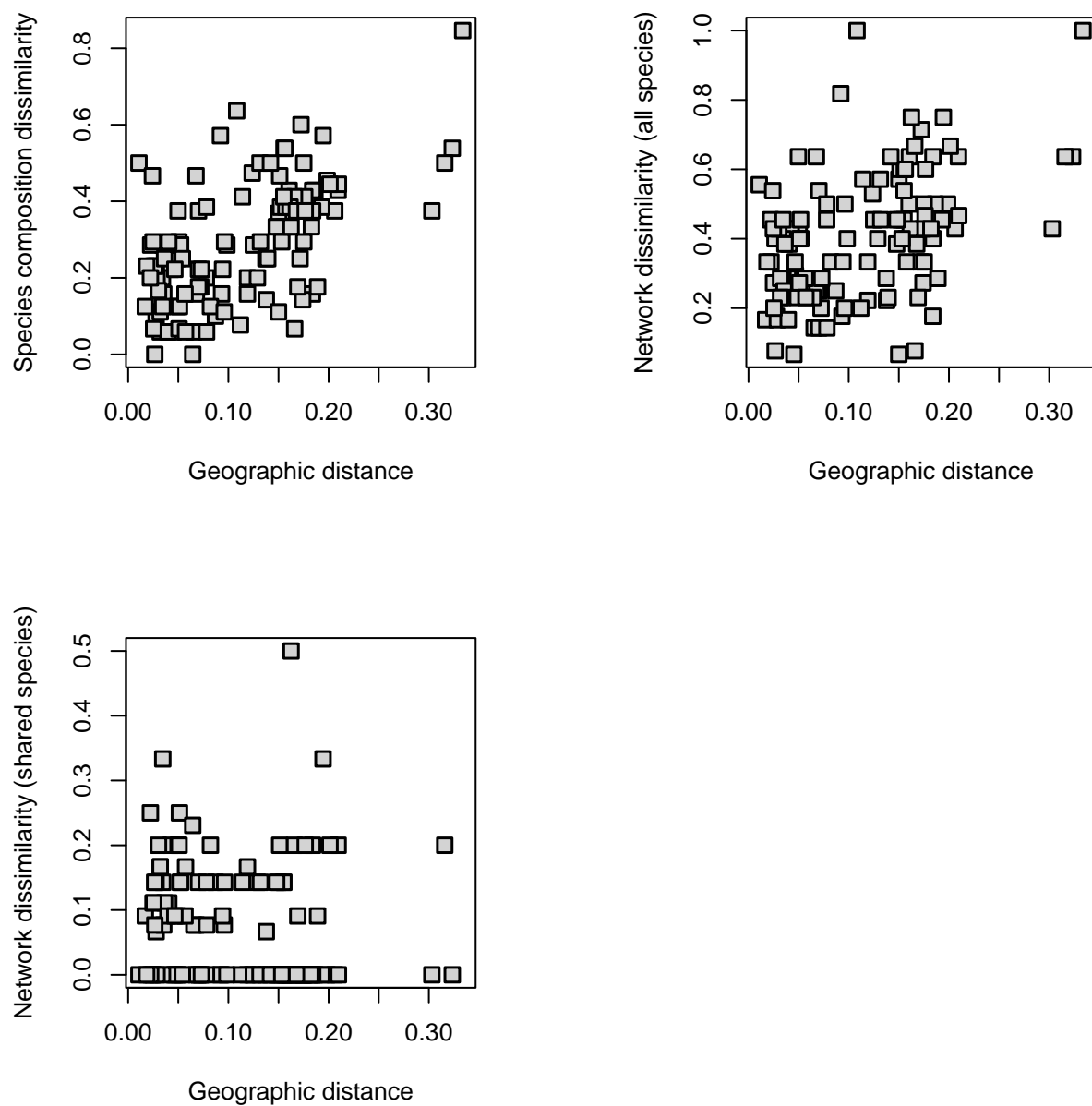


Fig. 3: Relationships between the geographic distance between two sites, and the species dissimilarity, network dissimilarity with all species, and network dissimilarity with only shared species.

- 1 over space is primarily driven by species turnover. The ease to gather both raw interaction data and associated meta-data
- 2 make conducting this analysis extremely straightforward.

3 Spatial visualization of networks

- 4 Bascompte (2009) uses an interesting visualization for spatial networks, in which each species is laid out on a map at the
- 5 center of mass of its distribution; interactions are then drawn between species to show how species distribution determines
- 6 biotic interactions. In this final use case, we propose to reproduce a similar figure (*Fig. 4*).

```
library(maps)
library(mapdata)
library(RColorBrewer)
library(sp)
library(plyr)
library(igraph)

# We fill a community data matrix
sp_by_site <- llply(graphs, function(x) unlist(V(x)$name))
sp_list <- unique(unlist(sp_by_site))
M <- matrix(0, ncol = length(sp_list), nrow = length(sp_by_site))
colnames(M) <- sp_list
rownames(M) <- names(sp_by_site)
for (site in c(1:length(sp_by_site))) M[names(sp_by_site)[site], sp_by_site[[site]]] = 1

# Next, we get the center position for each species
# (i.e. the mean position of the sites it occurs at)
sp_center <- adply(M, 2, function(x) colMeans(LatLon[names(x)[x > 0], ]))
rownames(sp_center) <- sp_center[, 1]
sp_center <- sp_center[, -1]

# We now create a regional network using betalink::metaweb
Mw <- metaweb(graphs)

# Plotting
source("suppmat/usecase_map.r")
```

7 Conclusions

- 8 The mangal data format will allow researchers to put together dataset with species interactions and rich meta-data, that are
- 9 needed to address emerging questions about the structure of ecological networks. We deployed an online database with
- 10 an associated API, relying on this data specification. Finally, we introduced rmangal, an R package designed to interact
- 11 with APIs using the mangal format. We expect that the data specification will evolve based on the needs and feedback

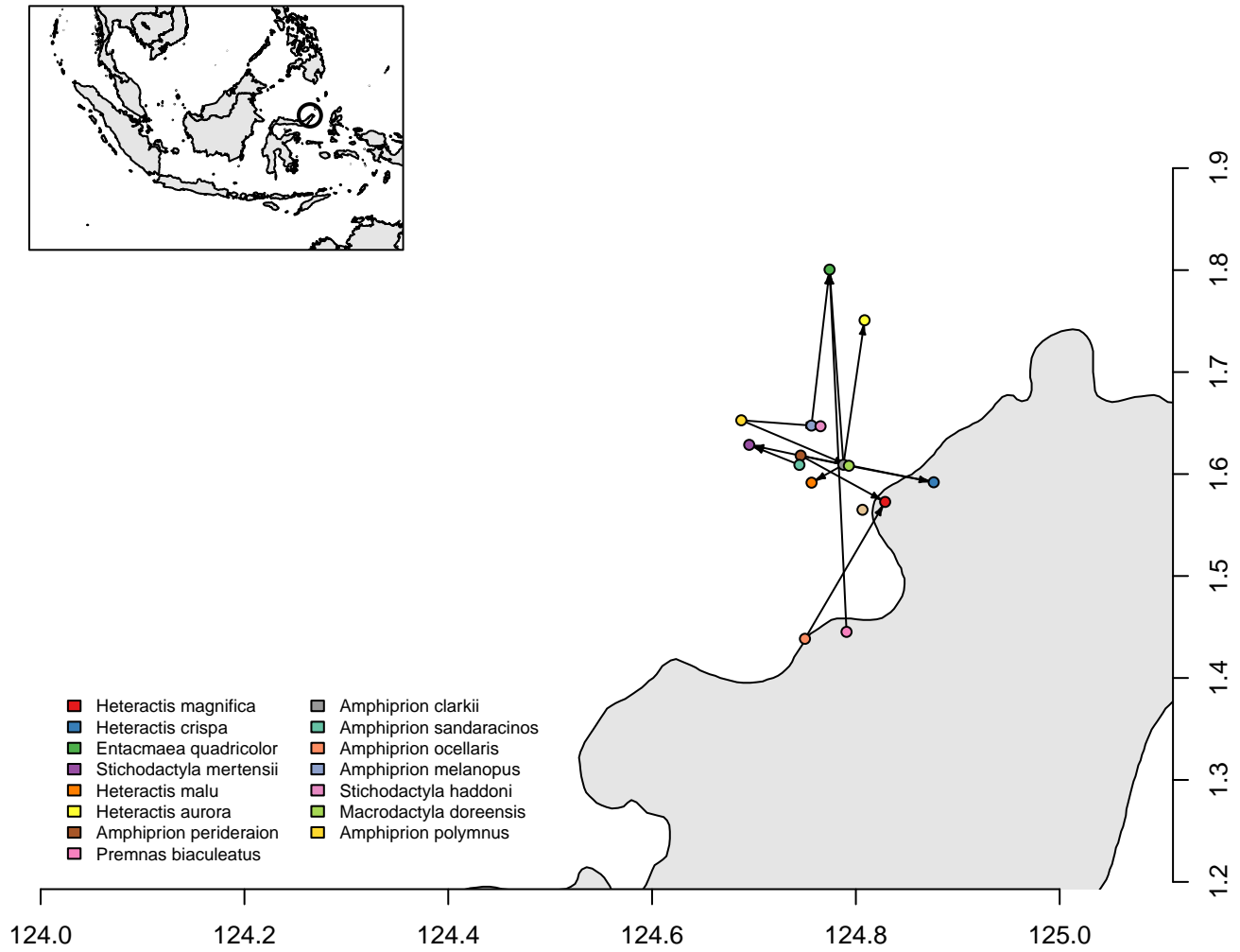


Fig. 4: Spatial plot of a network, using the maps and rmangal packages. The circles in the inset map show the location of the sites. Each dot in the main map represents a species, with symbiotic mutualisms drawn between them. The land is in grey.

1 of the community. At the moment, users are welcome to propose such changes on the project issue page: <https://github.com/mangal-wg/mangal-schemes/issues>. A python wrapper for the API is also available at <http://github.com/mangal-wg/pymangal/>. Additionally, there are plans to integrate this database with *GLOBI*, so that data
4 can be accessed from multiple sources (Poelen et al. 2014).

5 **Acknowledgements** This paper was developed during a workshop hosted at the *Santa Fe Institute*. TP, DBS, and DG
6 acknowledge funding from the Canadian Institute of Ecology and Evolution. We thank Scott Chamberlain and one
7 anonymous reviewer for comments on the manuscript. TP is funded by a start-up grant from the Université de Montréal.
8 We thank the rOpenSci team and developers for inspiration.

9 References

- 10 Bascompte, J. 2009. Disentangling the Web of Life. - *Science* 325: 416–419.
- 11 Bastolla, U. et al. 2009. The architecture of mutualistic networks minimizes competition and increases biodiversity. -
12 *Nature* 458: 1018–1020.
- 13 Chamberlain, S. A. and Szöcs, E. 2013. taxize: taxonomic search and retrieval in R. - F1000Research in press.
- 14 Chamberlain, S. A. et al. 2014. Traits and phylogenetic history contribute to network structure across Canadian plantpol-
15 linator communities. - *Oecologia*: 1–12.
- 16 Dalsgaard, B. et al. 2013. Historical climate-change influences modularity and nestedness of pollination networks. -
17 *Ecography* 36: 1331–1340.
- 18 Dunne, J. A. et al. 2002. Network structure and biodiversity loss in food webs: robustness increases with connectance. -
19 *Ecology Letters* 5: 558–567.
- 20 Eklöf, A. et al. 2013. The dimensionality of ecological networks. - *Ecology Letters* 16: 577–583.
- 21 Gravel, D. et al. 2013. Inferring food web structure from predatorprey body size relationships. - *Methods in Ecology and*
22 *Evolution* 4: 1083–1090.
- 23 Martinez, N. D. 1992. Constant connectance in community food webs. - *The American Naturalist* 139: 1208–1218.
- 24 Olito, C. and Fox, J. W. 2014. Species traits and abundances predict metrics of plantpollinator network structure, but not
25 pairwise interactions. - *Oikos*: n/a–n/a.
- 26 Pimm, S. L. et al. 1991. Food web patterns and their consequences. - *Nature* 350: 669–674.
- 27 Piwowar, H. A. and Vision, T. J. 2013. Data reuse and the open data citation advantage. - *PeerJ* 1: e175.
- 28 Poelen, J. H. et al. 2014. Global Biotic Interactions: An open infrastructure to share and analyze species-interaction
29 datasets. - *Ecological Informatics* in press.
- 30 Poisot, T. et al. 2012. The dissimilarity of species interaction networks. - *Ecology Letters* 15: 1353–1361.
- 31 Poisot, T. et al. 2013. Facultative and obligate parasite communities exhibit different network properties. - *Parasitology*
32 140: 1340–1345.
- 33 Poisot, T. et al. 2014. Beyond species: why ecological interaction networks vary through space and time. - *Oikos*:
34 n/a–n/a.
- 35 Reichman, O. J. et al. 2011. Challenges and opportunities of open data in ecology. - *Science* 331: 703–5.

- 1 Ricciardi, F. et al. 2010. Assemblage and interaction structure of the anemonefish-anemone mutualism across the Manado
2 region of Sulawesi, Indonesia. - *Environmental Biology of Fishes* 87: 333–347.
- 3 Rohr, R. P. et al. 2010. Modeling food webs: exploring unexplained structure using latent traits. - *The American naturalist*
4 176: 170–7.
- 5 Schleuning, M. et al. 2011. Specialization and interaction strength in a tropical plant-frugivore network differ among
6 forest strata. - *Ecology* 92: 26–36.
- 7 Tylianakis, J. M. et al. 2007. Habitat modification alters the structure of tropical hostparasitoid food webs. - *Nature* 445:
8 202–205.
- 9 Vines, T. H. et al. 2014. The Availability of Research Data Declines Rapidly with Article Age. - *Current Biology* 24:
10 94–97.