

1 **Title: A network approach for inferring species associations from co-occurrence data**

2 **Authors:**

3 Naia Morueta-Holme^{1,2,*} (morueta-holme@berkeley.edu)

4 Benjamin Blonder^{1,3} (bblonder@gmail.com)

5 Brody Sandel¹ (brody.sandel@bios.au.dk)

6 Brian J. McGill⁴ (mail@brianmcgill.org)

7 Robert K. Peet⁵ (peet@unc.edu)

8 Jeffrey E. Ott⁵ (jeott@live.unc.edu)

9 Cyrille Violle⁶ (cyrille.violle@cefe.cnrs.fr)

10 Brian J. Enquist⁷ (benquist@email.arizona.edu)

11 Peter M. Jørgensen⁸ (peter.jorgensen@mobot.org)

12 Jens-Christian Svenning¹ (svenning@bios.au.dk)

13

14 *Corresponding author.

15 NMH and BB contributed equally to this project.

16 **Affiliations:**

17 ¹Section for Ecoinformatics and Biodiversity, Department of Bioscience, Aarhus University, Ny Munkegade
18 114, DK-8000 Aarhus C, Denmark.

19 ²Integrative Biology, University of California – Berkeley, CA, USA.

20 ³Environmental Change Institute, School of Geography and the Environment, University of Oxford, Oxford,
21 UK.

22 ⁴School of Biology and Ecology/Sustainability Solutions Initiative, University of Main, Orono, ME, USA.

23 ⁵Department of Biology, University of North Carolina, Chapel Hill, NC 27599-3280, USA.

24 ⁶CEFE UMR 5175, CNRS - Université de Montpellier - Université Paul-Valéry Montpellier – EPHE -1919
25 route de Mende, F-34293 Montpellier, CEDEX 5, France.

26 ⁷Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, AZ, USA.

27 ⁸Missouri Botanical Garden, P.O. Box 299, St. Louis, MO 63166-0299, USA.

28 **Abstract**

29 Positive and negative associations between species are a key outcome of community assembly from regional
30 species pools. These associations are difficult to detect and can be caused by a range of processes such as
31 species interactions, local environmental constraints and dispersal. We integrate new ideas around species
32 distribution modeling, covariance matrix estimation, and network analysis to provide an approach to
33 inferring non-random species associations from local- and regional-scale occurrence data. Specifically, we
34 provide a novel framework for identifying species associations that overcomes three challenges: 1) finding
35 associations not driven by regional-scale distributions, 2) avoiding spurious associations caused by spatial
36 autocorrelation in regional null models, and 3) correcting for indirect effects from other associated species.
37 We highlight a range of research questions and analyses that this framework is able to address. The approach
38 is statistically robust using simulated data. In addition, we present an empirical analysis of >1,000 North
39 American tree communities that gives evidence for weak positive associations among small groups of
40 species. Finally, we discuss several possible extensions for identifying drivers of associations, predicting
41 community assembly, and better linking biogeography and community ecology.

42 **Introduction**

43 An unresolved question in ecology is how species assemblages come together at different spatial scales from
44 regional species pools to form local communities (Weiher et al. 2011, HilleRisLambers et al. 2012). Some
45 subsets of species may be consistently associated or dissociated in assemblages due to multiple processes
46 including chance, species interactions, and the indirect effect of showing the same (or opposite) response to
47 environmental conditions. While the effect of broad-scale environmental filtering and dispersal limitation
48 can be assessed using niche modeling techniques (De Marco et al. 2008, Guisan and Rahbek 2011, Normand
49 et al. 2011), the influence of species interactions and local environmental conditions on community assembly
50 is more challenging to measure and incorporate into predictions (Kissling et al. 2012, Pottier et al. 2013,
51 Wisz et al. 2013, Thuiller et al. 2013, Araújo and Rozenfeld 2014). If species interact with each other or
52 share resource or scenopoetic requirements (Soberón 2007) not adequately described by broad-scale models,
53 then stacking independent species distribution models to predict species assemblages (sensu Guisan and
54 Rahbek 2011, Calabrese et al. 2014) will provide misleading predictions of fine-scale community assembly.
55 Thus, a better understanding of species associations across scales could improve predictions of the dynamics
56 of local community composition in changing environments.

57 The goal of this paper is to improve the tools needed to detect interspecific associations from co-
58 occurrence data. We first briefly describe the development of co-occurrence methods and then draw from
59 different lines of research to present a more complete and flexible general framework for inferring species
60 associations that overcomes multiple challenges faced by previous approaches.

61 ***From experiments to co-occurrence methods***

62 Efforts to infer species associations and their role in structuring communities have a long history.
63 Traditionally, associations have been derived from small-scale field observations (e.g. MacArthur 1958,
64 Bullock et al. 2000) or manipulative experiments (e.g. Stewart and Aldrich 1951, Connell 1961). Such
65 methods may be suitable for studying associations among a few species at local scales. However, addressing
66 whether small-scale associations occur consistently across large regions raises major practical issues.

67 Experimental data is time-consuming to obtain even at small scales (e.g. Bullock et al. 2000, Callaway et al.
68 2002b), so these approaches are unfeasible for large numbers of species because the number of possible
69 interactions grows rapidly with assemblage size. For example, $n=100$ species have 4,950 possible pairwise
70 interactions and 161,700 possible three-way interactions.

71 An alternative approach for detecting associations is to gather occurrence data from field
72 observations and analyze species co-occurrence matrices by comparing the observed patterns to null models.
73 Such approaches were widely used during discussions of the “checkerboard” paradigm for forbidden
74 combinations of species on islands, with extensive debate over the need for null models to compare observed
75 association patterns to random expectation (Diamond 1975, Connor and Simberloff 1979). The recent
76 increase in availability of occurrence data has brought a renewed interest in using correlations and algorithms
77 to infer species associations and the mechanisms behind them (e.g. Bruelheide 2000, Blick and Burns 2009,
78 Blois et al. 2014). Since the 1970s, statistical approaches have been refined to simultaneously analyze co-
79 occurrence patterns of multiple species pairs (e.g. Gotelli and Ulrich 2010). The main idea is to create a
80 community matrix where rows are species, columns are sites and elements represent the observed
81 presence/absence or abundance of each species at each site. The matrix is then compared to a set of
82 randomized matrices in order to detect non-random co-occurrence patterns (Connor and Simberloff 1979,
83 Gotelli and Ulrich 2010). Modern updates to null-model-based co-occurrence approaches can test the effects
84 of environmental drivers, species interactions, or both in structuring communities; for approaches based on
85 randomized null models see e.g. Gotelli et al. 2010, Ulrich et al. 2012; for approaches based on analytical
86 null models see Araújo et al. 2011, Veech 2013.

87 *Three missing ideas*

88 One challenge that most previous co-occurrence approaches ignore is the potential effect of other species on
89 a particular pairwise association, i.e. indirect effects (Brown et al. 2004, Harris 2015). If two competing
90 species share a positive (or negative) relationship to a third species, their occurrences could be positively
91 correlated, and thus an indirect effect (correlation) is inferred when the true effect (partial correlation) is
92 actually negative (Brown et al. 2004, Schäfer and Strimmer 2005). A similar idea is implemented in ecology

93 for joint species distribution models (JSDMs; Pollock et al. 2014), where species associations are inferred
94 after accounting for the environment. However, JSDMs usually do not resolve indirect effects from other
95 species (but see the inversion approach of Harris 2015). The indirect association problem is well known in
96 other fields, such as association detection in large genomic and cell-signaling datasets. One solution is to use
97 Gaussian graphical models (Schäfer and Strimmer 2005) or modifications of them, which estimate the partial
98 correlations between e.g. each gene pair after taking into account the remaining genes (e.g. Dobra et al. 2004,
99 Friedman et al. 2008). The approach ensures estimation of conditional instead of joint associations and has
100 recently been extended to ecological association networks (Harris 2015).

101 A second challenge is the need for more robust null models. Species are not distributed randomly
102 across sites, but rather have regional geographic distributions that are constrained by climate and dispersal
103 limitations. In the past, most null models have been defined by simply resampling the observed species \times site
104 community matrix \mathbf{O} (e.g. Gotelli and Ulrich 2010, Borthagaray et al. 2014). However, this approach fails to
105 account for the additional broad-scale constraints, introducing unrealistic null expectations of spatial
106 independence across randomized sites that disregard spatial autocorrelation in species' distributions
107 (Legendre 1993, Lennon 2000). In the framework we present here, we suggest that alternative null models
108 can be defined to simulate a specific process. By incorporating regional structure that is contingent on e.g.
109 climate-based species distribution models, we can thus gain more confidence that associations are the
110 outcome of species interactions rather than shared environmental requirements. The integration of regional
111 and local community occurrence data with modern network methods has the potential for improving
112 predictive power in community ecology (see Baselga 2010, 2012 for a related approach).

113 A third missing idea is the incorporation of network theory. Network approaches have become
114 common in the study of protein interactions, social structure, etc. (Newman 2010), but have not been applied
115 widely to non-bipartite species association networks. The central idea is that any individual association
116 between species can be better understood in the context of the network of associations between all other
117 species. For example, species with many associations may be preferentially linked to those with few
118 associations, suggesting a scenario where some species act as hubs or keystones. Thus, the network of

119 species associations is a useful way to visualize the community as well as quantify changes in species and
120 multi-species patterns via node- and network-level statistics to answer more sophisticated questions about
121 associations.

122 Node-level statistics provide insight into species' roles. For example, the unweighted degree of each
123 species characterizes the number of its association partners. The ratio between the number of positive and
124 negative associations per species can provide a measure of a species' role in the network (e.g. as an attractive
125 "aggregator", if a species has more positive than negative associations, or a repulsive "segregator" of other
126 species, if otherwise). Network-level statistics also provide insight into the overall structure of the species
127 assemblage. For instance, modularity gives a measure of the overall structure of the network, indicating the
128 amount of division of the network into clusters of nodes that are densely connected to each other, but
129 sparsely connected to nodes in other clusters (Newman 2010, Borthagaray et al. 2014). Higher modularity
130 indicates that groups of species are more likely to be mutually associated. Finally, overall non-random
131 numbers of links assigned to each node in the network can be tested by comparing the degree distribution to
132 a binomial distribution, which in the limit of many species is equivalent to a chi-square test for deviation
133 from a Poisson distribution.

134 The approach we suggest thus builds on recent attempts to infer species associations from
135 occurrence data and network theory (e.g. Borthagaray et al. 2014). Here we provide a more flexible
136 framework that includes both positive and negative associations, and has the ability to test for deviations
137 from a range of regional-scale null models. The framework is useful for systems in which associations are
138 possible across any pair of the study species (i.e. not bipartite networks such as pollinator-plant interactions,
139 e.g. Bascompte 2010). We currently implement the framework only for symmetric associations (i.e. (+/+),
140 (0,0) or (-,-)). For asymmetric associations like predation (+/-) or commensalism (+/0), see Discussion.

141 *Interpreting associations*

142 The potential drivers of species associations can be better assessed using a network framework combined
143 with null modeling. This enables us to identify multi-species modules (i.e. non-random groups) of positively

144 or negatively associated species and assess the importance of particular species in shaping the modules,
145 ultimately improving predictions of community assembly. If we focus on plants, positive associations, or
146 aggregations of species, may be caused by biotic interactions such as facilitation between nurse trees and
147 seedlings, nitrogen fixation by certain species improving soil fertility for other species, or through shared
148 pollinators, seed dispersers, or facultative mutualisms with endophytic fungi (e.g., Afkhami et al. 2014).
149 Alternatively, positive associations may be indicative of shared local environmental requirements or effects
150 of stabilizing niche differences (Chesson 2000, Lasky et al. 2014), or reflect historical dispersal dynamics
151 such as the expansion from glacial refugia (Svenning and Skov 2007). Negative associations, or segregation
152 of species, may also be driven by biotic interactions through competition, or alternatively reflect different
153 local requirements (e.g. variation in microclimate or edaphic conditions).

154 Once we have inferred the patterns of species associations, we can proceed to determine which
155 hypothesized underlying drivers are most important (Box 2), with the aim of understanding the principles
156 that cause associations to occur. Functional trait predictors may help identify mechanisms. For instance, if
157 positive associations reflect complementary niches where each species occupies a different niche in relation
158 to resource use (Bolnick et al. 2011), then phylogenetic (Webb et al. 2002) or functional (Weiher et al. 2011)
159 distances should be largest between positively associated species (Violle et al. 2011). Likewise, by inferring
160 spatial patterns of community-weighted mean associations (Box 2) we can test whether associations are
161 influenced by climate gradients, and are e.g. more common in warmer areas (Brown et al. 1996, Schleuning
162 et al. 2012), or whether positive associations become more prevalent at higher elevations (Callaway et al.
163 2002a).

164 **A framework for detecting associations**

165 Drawing on the multiple lines of research described above, we here propose a general framework to uncover
166 species associations from co-occurrence patterns. The associations we identify are those not explained by
167 broad-scale climate gradients. The main idea of our approach is to pair broad- and local-scale co-occurrence
168 information with the Gaussian graphical model approach, spatially explicit regional null models (Guisan and

169 Zimmermann 2000, Gotelli and Ulrich 2010, Borthagaray et al. 2014) and network theory (Newman 2010).
170 Our framework is implemented as the ‘netassoc’ R package.

171 We first show how the co-occurrence framework can infer a network of species associations by
172 pairing regional occurrence data with local-scale assemblage data. Second, we perform a simulation analysis
173 demonstrating that the framework has acceptable error rates in realistic use cases. Third, we demonstrate a
174 range of network-based analyses that can describe associations and test hypothesized mechanisms. Lastly,
175 we illustrate the framework with an application to the trees of eastern United States using a large dataset of
176 local co-occurrences and regional occurrences.

177 In Box 1 we describe the main data inputs and methodological choices for the approach. Briefly, the
178 first step is to identify non-random species associations for an observed and a null dataset. To do this, we
179 compare the partial correlation coefficients inferred for observed local-scale data to those inferred for
180 regional-scale null expectations calculated from independent data. We then reinterpret these co-occurrence
181 effect sizes within a network framework.

182 *Network construction*

183 To compute the species association network we need (i) the observed co-occurrences, i.e. a set of observed
184 presence-absence or abundance data for n species found in a random sample of m sites, and (ii) the expected
185 co-occurrences for the same species and sites based on a null model. The null expectation can be for
186 presence or abundance of the species. From (i) we generate \mathbf{O} , the observed species \times site community
187 matrix. From (ii) we generate \mathbf{N} , the expected abundance (or presence) patterns at each local site as predicted
188 by a chosen regional species distribution model (iii). Both the \mathbf{O} and \mathbf{N} matrix will have n rows (species) and
189 m columns (sites).

190 We first infer the association strength between species i and j as entries A_{ij} in the $n \times n$ matrix \mathbf{A} . We
191 calculate an $n \times n$ covariance matrix $\mathbf{\Sigma}$ for each of \mathbf{O} and \mathbf{N} . From the inverse of this covariance matrix we
192 obtain standard partial correlations between species, i.e. as

193 (1)
$$C_{ij}(\mathbf{M}) = \frac{-\Sigma_{ij}^{-1}(\mathbf{M})}{\sqrt{\Sigma_{ii}^{-1}(\mathbf{M}) \cdot \Sigma_{jj}^{-1}(\mathbf{M})}}$$

194 where \mathbf{M} is the $n \times m$ input community matrix. C_{ij} represents the effect of species i on species j after
 195 correcting for the effects of all other species; it is zero if species i is conditionally independent of species j .
 196 Equation 1 represents the fundamental mathematical approach taken when constructing Gaussian graphical
 197 models (Schäfer and Strimmer 2005) for inferring linear associations between random variables. We
 198 calculate $C(\mathbf{O})$ as well as for $K \gg 1$ resamples of \mathbf{N} , $C(\tilde{\mathbf{N}})$. This distribution simply simulates a weighted
 199 lottery model of community assembly where species enter a community based only on their overall
 200 abundance in the regional pool (i.e. the probabilities in \mathbf{N}). To do so, the resamples preserve the total number
 201 of individuals within each site, weighting the sampling by the expected abundance of each species based on
 202 the original \mathbf{N} matrix.

203 We then determine if the observed association between each species pair is positive or negative by
 204 comparing the observed co-occurrence score to the distribution of expected co-occurrences across $\tilde{\mathbf{N}}$. We
 205 calculate a standard effect size $SES_{ij}(\mathbf{O}, \mathbf{N})$, i.e., by rescaling by the mean (μ) and standard deviation (σ) of the
 206 null distributions:

207 (2)
$$SES_{ij}(\mathbf{O}, \mathbf{N}) = \frac{C_{ij}(\mathbf{O}) - \mu[C_{ij}(\tilde{\mathbf{N}})]}{\sigma[C_{ij}(\tilde{\mathbf{N}})]}$$

208 Finally, we distinguish between significant and non-significant associations. We first calculate a two-tailed
 209 p -value for species i and j as the number of replicates in which the absolute observed association strength is
 210 smaller than the absolute null association strength divided by the total number of replicates. The next step is
 211 to then correct for multiple comparisons by specifying a false discovery rate, α , and performing a Benjamini-
 212 Hochberg correction (Benjamini and Hochberg 1995) on each p -value, producing a new set of p -values p_{ij}^* .
 213 The false discovery rate conceptualizes the type I error rate by controlling the expected proportion of false
 214 positives, i.e. the proportion of incorrect rejections of the null hypotheses across multiple comparisons.

215 Together, this process yields a species-by-species adjacency matrix \mathbf{A} with n rows and n columns
 216 (number of species):

$$217 \quad A_{ij} = \begin{cases} p_{ij}^* < a & \text{SES}_{ij}(\mathbf{O}, \mathbf{N}) \\ p_{ij}^* \geq a & 0 \end{cases} \quad (3)$$

218 This matrix A is treated as the adjacency matrix (showing which species are connected to each other) used to
 219 define the species association network, such that a significantly positive or negative association between
 220 species i and j is established with strength $A_{i,j}$ if $A_{i,j}$ is nonzero. The network of associations is used for all
 221 subsequent analyses.

222 *A few important decision points in the framework*

223 While our statistical framework for describing species associations is general, the user must choose between
 224 multiple definitions and parameters specific to the system and taxa being studied. For example, one
 225 important choice in the analysis is the type of null model. The \mathbf{N} matrix shows the predicted abundance or
 226 presence/absence at each local site for each species. Multiple methods can be used to define \mathbf{N} . For example,
 227 a leave-one-out LOESS model on occurrence data can be used to calculate the expected abundance at one
 228 site from a distance interpolation of the observed abundances at all other sites. Such a model indicates the
 229 expected community produced by a dispersal-environment model. Alternatively, MaxEnt or other species
 230 distribution models can be used to calculate the expected abundance based on only broad-scale climate.
 231 Other approaches are possible, like calculating the expected occurrence from stacking species' regional
 232 geographic ranges sourced from expert-drawn range maps or from mechanistic regional models that predict
 233 abundance patterns across space (e.g. demographic or trait-based dispersal models; Jongejans et al. 2008).

234 A second set of choices that the user must define in the analysis is how to estimate the inverse
 235 covariance matrix Σ^{-1} , which can be difficult in practice. Two situations can arise: first, the number of
 236 species can be much larger than the number of sites; second, most sites can contain very few species. Both
 237 cases can lead to Σ becoming singular (i.e. non-invertible) because of very large covariances between some

238 species pairs. A range of shrinkage estimators for Σ^{-1} have been developed that provide a robust approach to
239 resolving this general problem (e.g., Hoerl and Kennard 1970, Schäfer and Strimmer 2005, Friedman et al.
240 2008). All the shrinkage estimators increase estimator bias in exchange for reduced mean squared error by
241 introducing additional offset parameters that ‘shrink’ coefficient estimates and so force the existence of a
242 matrix inverse. These offset parameters can be estimated by cross-validation approaches. A full survey of
243 these methods is beyond the scope of this article, but some popular options include the James-Stein type
244 shrinkage estimator (Schäfer and Strimmer 2005), the graphical lasso (L_1 -regularization; Friedman et al.
245 2008), or ridge regression (L_2 -regularization; Hoerl and Kennard 1970). We caution against using the
246 graphical lasso because it produces sparse inverse covariance matrices that can produce singular null partial
247 correlation distributions $C_{ij}(\tilde{\mathbf{N}})$, and instead recommend using the James-Stein estimator because of its good
248 performance and low computational cost (Schäfer and Strimmer 2005).

249 Finally, we also recommend log-transforming abundance data as $f(x, a) : x \mapsto \log(x + a) - \log(a)$
250 for some small number a , e.g. 10^{-6} . This transformation can improve normality of the distribution of
251 abundance data, which can otherwise take either zero or very large values.

252 **Testing the framework with a simulation analysis**

253 To measure the expected performance of the network framework, we simulated co-occurrence matrices with
254 known associations and determined how well network-detected associations matched these.

255 Consider a scenario involving n species distributed across m sites, of which a fraction h are
256 unsuitable. We first generated an $m \times n$ expected species-by-site matrix \mathbf{N} , all of whose abundance entries
257 were independently and identically distributed according to a hurdle model, such that N_{ij} was zero with
258 probability h and Poisson-distributed (mean λ) with probability $1-h$. If any of the marginal sums of \mathbf{N} were
259 zero (i.e. a site with no species or a species with no sites; problematic only for small n and m), we re-
260 generated \mathbf{N} until all marginal sums were non-zero. We then independently generated an $m \times n$ observed
261 matrix \mathbf{O} via the same hurdle process.

262 As a next step, we assumed that there were Z associations in the ‘true’ association network. We
 263 chose the Z associations by generating a random graph with n vertices and Z edges (Erdős and Rényi 1959),
 264 with weight w_z ($z \in \{1,2,\dots,Z\}$) set to either 1 or -1 with equal probability. To model this, we iterated over all
 265 associations z ; for each pair of species i_z and j_z for which a true association exists, we chose a random
 266 fraction f of sites $\{m_z\}$; at each of these sites we either increased the abundance for both species (when $w_z>0$)
 267 or increased for one species and decreased for the other (when $w_z<0$) by a factor s proportional to the mean
 268 abundance of both species at these sites:

$$\begin{aligned}
 \mathbf{O}_{i_z, m_{\{z\}}} &= \mathbf{O}_{i_z, \{z\}} + w_z^2 \cdot s \cdot (\mathbf{O}_{i_z, \{z\}} + \mathbf{O}_{j_z, \{z\}}) / 2 \\
 \mathbf{O}_{j_z, m_{\{z\}}} &= \mathbf{O}_{j_z, \{z\}} + w_z \cdot s \cdot (\mathbf{O}_{i_z, \{z\}} + \mathbf{O}_{j_z, \{z\}}) / 2
 \end{aligned}$$

269 (5)

270 This process effectively increased the covariance between species when the species were positively
 271 associated and decreased it when the species were negatively associated, with the parameters h and s
 272 controlling the strength of the association. We used $h=0.2$ and $s=0.2$ in this analysis.

273 Next, we applied our network framework using the matrices for all possible parameter combinations
 274 of $n=10, 100$; $m=10, 100, 1,000$, $f=0, 0.5$, and $Z=10, 50, 100$ (note that some combinations were not possible,
 275 e.g. when $n=10$, we cannot simulate values of $Z>45$ because they would exceed the number of interactions in
 276 a fully connected network). We used a James-Stein type shrinkage estimator (Schäfer and Strimmer 2005)
 277 with significant associations inferred at the $\alpha=0.05$ significance level. We set the number of null replicates to
 278 1,000 and repeated the entire analysis for each parameter combination 10 times.

279 In order to calculate error rates for our method, we compared the inferred networks’ structure to the
 280 true network’s structure. We counted a true positive association if it was detected for the correct pair of
 281 species and had the correct sign; as a true negative if it was not detected for a pair of species for which an
 282 association did not exist. A false positive association was counted if it was detected but was either the
 283 incorrect sign or the pair of species did not have a true association. Similarly, a false negative was counted if
 284 it was not detected but the pair of species did have a true association. These counts allowed us to calculate

285 the positive predictive value (PPV; true positives divided by true positives plus false positives) and the
286 negative predictive value (NPV; true negatives divided by true negatives plus false negatives) as summary
287 statistics.

288 To determine the sensitivity of the method to different parameters, we constructed a random forest
289 regression model for NPV and PPV. The method generates an ensemble of regression trees and therefore
290 allows for multi-way interactions between variables. We calculated the importance of each variable as the
291 residual sum of squares caused by splitting on the variable of interest, averaged over all trees. Random forest
292 models were built using the *randomForest* R package (Liaw and Wiener 2002).

293 Across all parameter combinations, PPV took a mean value of 12 ± 26 s.d. %, while NPV took a
294 mean value of 87 ± 26 s.d. %. That is, the method was better at detecting the absence of associations than the
295 presence of associations. A random forest model predicting PPV as a function of m , n , f , and Z explained
296 56% of the variation in the data and showed that m had the largest impact (14.6) on PPV, followed by Z
297 (8.1), with smaller contributions of n (1.7) and f (1.3). Partial dependence plots indicated that larger values of
298 each parameter led to higher values of PPV. A similar random forest model for NPV explained 53% of the
299 variation in the data, and showed that m (11.1) and Z (8.7) had the largest impacts, with smaller contributions
300 of n (1.5) and f (2.6). Partial dependence plots indicated that smaller values of each parameter led to higher
301 values of NPV. Additionally, datasets with large fractions of zero-abundance records did not challenge the
302 model's ability to infer associations.

303 Overall, this analysis indicates that the method trades off between successfully detecting true
304 associations (high PPV) and successfully detecting the absence of false associations (high NPV). Better PPV
305 occurs for large datasets, while better NPV occurs for small datasets. Increasing the value of the false
306 discovery rate α can further control this tradeoff. We repeated analyses for $\alpha=0.5$ (results not shown), which
307 approximately doubled PPV and halved NPV in all cases. Thus, the method can yield acceptable
308 performance in a wide range of realistic use cases.

309 **Empirical test of the framework**

310 After establishing the robustness of the framework through the simulations, we can apply it to real datasets
311 and analyze the inferred associations. Here we present an illustration of the approach using communities
312 from temperate forests in North America. In particular, we use the framework to test whether i) there are
313 positive and/or negative associations among tree species that cannot be explained either from a regional
314 dispersal-environmental model, or from broad-scale climate gradients; and whether ii) such associations can
315 be explained by phylogenetic relatedness, functional trait similarity, or environmental gradients. We predict
316 that across null models, association networks will have non-random structure and high modularity, and that
317 positive (negative) associations will be more common among more (less) closely related or more similar
318 species or in warmer (colder) environments.

319 *Species data*

320 We chose local community data to represent eastern North America. We extracted community-level tree
321 species abundance data from the Forest Inventory and Analysis (FIA) database (Gray et al. 2012) (see
322 Appendix S1 for query details). We selected 5,138 0.07 ha plots from the eastern USA (east of the 100° W
323 meridian), and subsampled to a maximum of three plots for each 10,000 km² to reduce spatial sampling
324 biases. Each plot consists of four 7.3 m radius subplots located 36.6 m from each other. All plots included
325 were surveyed in the field between 2004 and 2008, used standard sampling protocols, and were marked as
326 natural stands without evidence of artificial regeneration or human disturbance. If a plot was surveyed more
327 than once in the time period chosen, we only included the newest survey. We excluded FIA taxa that were
328 not identified to the species level. A total of $m=1,009$ plots out of the 5,138 plots and $n=137$ tree species
329 were included in the analyses.

330 Point occurrence data for the null models came from BIEN, the Botanical Information and Ecology
331 Network (Enquist et al. 2009, <http://bien.nceas.ucsb.edu/bien/>) for each of the 137 tree species.

332 *Alternative species distribution models*

333 We computed results based on three definitions of the null model. First, we used the commonly used
334 random-swap algorithm, where the local community matrix \mathbf{O} is randomized to create the null matrix \mathbf{N}

335 (e.g., Connor and Simberloff 1979, Gotelli and Ulrich 2010), keeping row and column sums fixed (i.e. total
336 species and site abundances). Second, we used a leave-one-out LOESS regression of plot abundances with a
337 span parameter of 0.2 (mirroring ter Steege et al. 2013) to compute the expected abundance at each plot from
338 its spatial position and those of the remaining 5,137 plots in the full FIA dataset. Third, we also used species
339 distribution models created with the algorithm MaxEnt (Phillips and Dudík 2008) and based on the BIEN
340 point occurrence data to estimate the climatic potential range of each species. As a test case, we used 19
341 bioclimatic layers representing “current” climate (average 1950–2000 conditions) as model predictors,
342 extracted from WorldClim 1.4 at 30-arc second resolution (Hijmans et al. 2005). Each model was fit using
343 default parameters to both North and South America to capture the climatic range of the full New World
344 distribution of each species. We converted the model suitability scores to expected abundance values by
345 standardizing them so that the summed suitability scores for each species equaled the total number of
346 individuals across all plots. This procedure assumes a linear relationship between suitability and abundance.

347 We used 1,000 resamples \tilde{N} of the expected species \times site matrix N . As in the simulation analysis,
348 we used a James-Stein shrinkage estimator for the inverse covariance matrix, and specified an overall false
349 discovery rate of $\alpha=0.05$ to exclude non-significant associations. Modules in the network were inferred using
350 a standard fast-greedy algorithm (Clauset et al. 2004).

351 *Potential predictors of species' associations*

352 We obtained data on phylogenetic relatedness between species and functional trait values to illustrate how
353 the network statistics can be used to test hypotheses of drivers of species associations.

354 We calculated phylogenetic distances from a phylogeny for all the trees of eastern USA using
355 Phylocom's 'phylomatic' tool (Webb et al. 2008). We used the R20120829 backbone tree, with branch
356 lengths adjusted by fossil constraints (Gastauer and Meira-Neto 2013). We then computed the distance
357 between each pair of species.

358 We also obtained measurements of four traits thought to underlie major axes of ecological strategy
359 variation (Westoby et al. 2002): maximum height (m), specific leaf area (SLA; cm^2/g), seed mass (g) and

360 wood density (g/cm^3). Functional trait data were extracted from the BIEN database. There was good
361 coverage for trait data: 99% for height, 72% for SLA, 93% for seed mass, and 81% for wood density. We
362 \log_{10} -transformed each trait value to reduce skewness, then rescaled all values by subtracting means and
363 dividing by standard deviations. We then computed trait differences between all pairs of species.

364 We used linear regression to determine whether the links between each species pair (link strength if
365 adjacent; 0 if not adjacent) was predicted by pairwise phylogenetic or trait distance between the species pair.
366 We did not correct for non-independence of predictor distances (e.g. Mantel-type test). This approach should
367 lead to an increased rate of falsely rejecting the null hypothesis, meaning that failing to reject the null
368 hypothesis is more likely to reflect a true absence of relationship.

369 *Network structure of trees of eastern USA*

370 The species association network based on the random-swap algorithm showed a non-random structure (chi-
371 square test for Poisson degree distribution, $p < 10^{-45}$), identifying many positive and negative associations
372 (Table 1). On the other hand, the overall structure of networks based on the two regional null models – the
373 LOESS dispersal-environment model, and the MaxEnt climate model – was not different from random. This
374 result indicates that the overall distribution of local associations between trees in North America can mainly
375 be explained by broad-scale drivers. Looking at the individual associations, we only found a few positive
376 interactions in both of these networks. Interestingly, all five modules of 2-3 associations identified using the
377 LOESS regional model also appeared when applying the MaxEnt model, although sometimes with an
378 additional species in the module (Fig. S1). We found additional associations deviating from the MaxEnt
379 regional model, and modules were larger in general (Table 1, Fig. 1). In the random-swap network most
380 species were inferred to be ‘segregators’, while in the MaxEnt and LOESS network most species were
381 inferred to be ‘aggregators’ (Fig. S2).

382 *What predicts species’ associations?*

383 The positive and negative associations we found were not predicted by phylogenetic or trait distances (Fig.
384 S3) regardless of null model. The network structure explained by phylogeny in any network was no more
385 than 0.05%, and the variation explained by all four traits together was no more than 0.12% in any network.

386 We did find weak spatial gradients in abundance-weighted mean values of degree for the species
387 comprising each community (Fig. S4). For each of the random-swap and MaxEnt networks, mean degree
388 was negatively correlated with mean annual temperature but not with mean annual precipitation (multiple
389 regression; both $p < 0.002$, both $R^2 < 0.06$). The LOESS network was not correlated with either mean annual
390 temperature or precipitation ($p=0.26$).

391 **Can we infer non-random species associations?**

392 Our simulation analysis showed that the co-occurrence framework is indeed able to identify known
393 associations within the parameter regimes we explored. However, there is a trade-off, where large datasets
394 lead to better success in detecting true associations, and smaller datasets are better at identifying the
395 absence of false associations. When including explicit regional null models, the framework is able to
396 pinpoint which species are associated after broad-scale drivers have been accounted for. This not only allows
397 testing for multiple regional models. Another advantage is that spatial autocorrelation in species distributions
398 can be taken into consideration. Indeed, when looking at the results from our case study, it is apparent that
399 the random-swap algorithm identifies spurious positive and negative associations that stem from ignoring the
400 spatial dependence of species occurrences. This important implication shows that the null model from a
401 hypothesized process such as dispersal or broad-scale climate patterns gives better control of the type of
402 associations identified.

403 Another challenge addressed by the framework is that of indirect effects from multiple species on
404 pairwise associations. The result is that when applying the framework to tree communities of eastern North
405 America, we only found a small number of positive associations deviating from either the LOESS or the
406 MaxEnt based regional models after correcting for indirect effects (Fig. 1). Ecologically, this may reflect that
407 tree species distributions are largely controlled by environment and dispersal, with little importance of

408 interspecific interactions, even if the latter matter for local abundances. The associations we do identify are
409 those deviating from the broad-scale expectations. The LOESS model can be interpreted as a dispersal and
410 environment model, since it simply interpolates abundance as a function of distance and implicitly includes
411 environmental conditions that are spatially autocorrelated. Alternatively, the MaxEnt model computes the
412 expected co-occurrence as a function of broad-scale climate variables, ignoring dispersal constraints, and
413 instead representing the potential climate range of each species. It is thus not entirely surprising that the
414 associations identified with the dispersal-environment model are all a subset of those identified with the
415 climate-only model. Contrary to our predictions, functional traits, phylogenetic relatedness and
416 environmental gradients did not correlate with the associations found across the networks. However, natural
417 history, successional dynamics and missed environmental drivers could explain at least some of the
418 associations identified. Indeed, local habitat requirements not fully captured by the broad scale
419 environmental gradients tested here seem to explain several associations. For instance, *Carya aquatica* and
420 *Taxodium distichum* are both species of Coastal Plain, strongly associated with large river backswamps that
421 are periodically flooded. A similar habitat is preferred by *Nyssa aquatica* and *Planera aquatica*, although
422 this species pair is found more up-stream than the first ones, in areas where flooding is less prolonged.
423 *Quercus palustris*, *Q. bicolor* and *Carya laciniosa* also prefer swamp habitats, though more from interior
424 flatlands on more calcareous soils. Habitat preference does not seem to explain the association between
425 *Carya pallida* and *Quercus michauxii*, which prefer dry/sandy and swampy soils, respectively. These two
426 species rarely occur together, but both are largely confined to the southeastern Coastal Plain and lower
427 Piedmont regions. This geographic signal could instead be driven by dispersal limitation. A few boreal
428 modules such as the one around *Abies balsamea* are identified with the MaxEnt regional model but disappear
429 in the LOESS model. These modules may represent local mosaics of boreal and temperate stands at the
430 boreal-temperate transition zone, reflecting local environmental conditions and/or priority effects (Pastor and
431 Mladenoff 1992).

432 Given that the MaxEnt model only covers environmental constraints, species might be simulated to
433 co-occur that have similar climatic requirements but are allopatric. In such a case we would expect more

434 negatively associated species in the network derived from the MaxEnt null model. Instead, the modularity of
435 networks based on the MaxEnt models was much higher than that of LOESS-based networks, and in neither
436 case did we find negative associations. The fact that we only identify positively associated groups of species
437 unexplained by broad-scale climate conditions is consistent with the results of the simulation study of Araújo
438 and Rozenfeld (2014), who found that the effect of positive (but not negative) dependencies between species
439 scaled up to biogeographical scales and should be accounted for in range models under climate change.

440 One implication of our results – in particular the larger amount of positive associations deviating
441 from the climate-only model – is that we cannot rely solely on stacked species distribution models (SDMs)
442 (sensu Guisan and Rahbek 2011, Calabrese et al. 2014) to predict the composition of local communities, but
443 need to take species associations into account (Araújo and Rozenfeld 2014) – whether driven by biotic
444 interactions, dispersal, or local environmental filtering. In cases where e.g. local environmental or dispersal
445 data are available, these can be integrated to filter broad-scale predictions and improve the performance of
446 SDMs for communities (Boulangéat et al. 2012).

447 We note that even low error rates can translate into high absolute numbers of incorrect associations
448 for datasets with large species numbers. Given that 100 species have 4,950 potential pairwise interactions,
449 our simulation error rates mean that, on average, approximately 25 links will be inferred that do not actually
450 exist. Our Type 2 error rates are very high for small datasets, but settle to near zero for datasets with at least
451 100 sites, with performance further increasing when species co-occurrence patterns have high association
452 weights. Thus, the algorithm does not often miss real associations, but even then, our Type 2 error rates
453 translate into approximately 500 real links that are missed. These rates and uncertainties are similar to other
454 association-inference approaches (Morales-Castilla et al. 2015) and suggest that these frameworks are most
455 useful for reducing the space of possible associations to significantly more manageable numbers (cf. Figure 2
456 in Morales-Castilla et al. 2015).

457 *Inferring drivers of associations — a complex challenge*

458 In our example with North American tree communities, the associations we identified can best be explained
459 by ecological (and possibly geographic) groupings. It is thus not surprising that they could not be predicted
460 by markers like functional trait similarity and phylogenetic distance. Similarly, we only found a weak
461 negative correlation between community mean degree and mean annual temperature for the associations
462 identified with the MaxEnt null model, possibly reflecting the geographic signal of boreal communities. The
463 application still shows how network metrics can be used to test for drivers of associations, although in this
464 particular case, more direct measures such as species' flooding tolerance from Ellenberg values would have
465 been more useful.

466 Still, we can imagine several issues that could pose obstacles to the prediction of associations. One
467 issue is that multiple competing processes can be at play. Competition and facilitation may cancel each other
468 out (Callaway and Pennings 2000). Species associations may also vary across environmental gradients
469 (Callaway et al. 2002a, Pottier et al. 2013) or across temporal scales (Martorell and Freckleton 2014, Blois et
470 al. 2014), such that the scale of input data used would be critical. Even in networks of biotic interactions
471 between plants and pollinators, pairwise associations do not always correlate with hypothesized drivers, even
472 when correlations are found with metrics of overall network structure (Olito and Fox 2015). We therefore
473 expect that disentangling the processes driving patterns of co-occurrence will remain an ongoing challenge.

474 *Extensions of the network framework*

475 The increased availability of regional species occurrence and local-scale co-occurrence data across
476 large extents has been a strong driver behind recent attempts to disentangle the processes promoting species'
477 associations through biotic interactions, dispersal limitations and environmental filtering (e.g. Blick and
478 Burns 2009, Ulrich et al. 2012, Blois et al. 2014). We have presented a flexible framework to infer species
479 positive and negative associations that deviate from expected broad-scale processes, and illustrated how
480 network statistics can be used to test hypothesized drivers of associations. The approach can be readily
481 applied to other datasets and systems for any other type of potential species associations across trophic
482 levels.

483 There is an ongoing debate on the directionality of interactions, with some studies finding more
484 positive than negative associations (Blick and Burns 2009 and references therein, Blois et al. 2014), while
485 others (including a meta-analysis) have found the opposite (Azaele et al. 2010, Gotelli and Ulrich 2010).
486 Inconsistent associations appear to be the norm rather than the exception in the literature, changing when
487 using species-based rather than individual-based models (Blick and Burns 2009), or when looking for
488 patterns across time in the paleo-record (Blois et al. 2014). Some of the disparate results seen across studies
489 might be the unintended outcome of methodological differences. We suggest that sensitivity analyses and
490 consensus approaches should be used to achieve robust results. Importantly, partial correlations should be
491 widely implemented, as well as taking into account spatial autocorrelation as part of the null model to avoid
492 identifying non-existent interactions.

493 The framework as currently implemented only accounts for symmetric associations, e.g. those
494 representing (+/+), (0,0) or (-,-) interactions. This limitation is imposed because the matrix **A** is derived from
495 **C**, which is derived from Σ^{-1} , which is symmetric. The coefficients in **C** could instead be calculated using
496 other approaches that allow for non-symmetric outcomes, and so also capture effects like predation (+/-) and
497 commensalism (+/0). The network framework allows for such directed linkages, and the R package does
498 allow arbitrary user-specified functions to be used for calculating **C**. However, we are unaware of viable
499 candidate functions to use. Most association metrics are symmetric (cf. Janson and Vegelius 1981) and the
500 few that are not (e.g. Somers' D (Somers 1962), Goodman and Kruskal's lambda (Goodman and Kruskal
501 1972)) cannot be used with abundance data, do not take both positive and negative values, and do not
502 account for indirect associations. Aside from the excess co-occurrence approach of Araújo and Rozenfeld
503 (2011), the Markov network approach suggested by Harris (2015), and the directed partial correlation
504 coefficients proposed for microarray studies (Yuan et al. 2011), we are unaware of other approaches.
505 Developing association metrics that capture all possible ecological integrations types should be a priority.

506 For cases in which many strong associations are identified, the network framework could potentially
507 be extended to provide predictions of local community composition, which in turn can be tested in new
508 communities. For instance, with information on the identity of only one species in a community, a “network-

509 crawling” approach could be used to predict the identity of the remaining species by allowing species that are
510 close to the one known species in the network are more likely to be predicted in the local community (or less
511 likely for species negatively associated). Such an extension would provide transparent and powerful tests of
512 the predictive ability of association networks derived either solely from occurrence data, or combined with
513 experimental or field-based interaction data.

514 More sophisticated analyses could be applied to our framework to provide a more thorough test of
515 the drivers of associations observed. Indeed, a recent study found modularity analyses useful to identify
516 biological attributes driving the connection of modules of co-occurring species (Borthagaray et al. 2014).
517 Such additional analyses can be readily applied with our framework, potentially combined with other
518 methods to identify network modules (Leger et al. 2015) to e.g. investigate the relationship of traits and
519 phylogenetic relationships for species within each module, and thus explore the definitions of the scale of
520 study.

521 *Looking forward*

522 Mechanistic understandings of the drivers of species distributions across scales are needed to better
523 predict the consequences of rapid global change. The general framework we propose here provides a novel
524 tool to infer local-scale associations that can affect species’ distributions. We have presented a flexible
525 framework linking co-occurrence and null-models to network theory for inferring species associations that is
526 easily applicable to other systems and datasets. This approach can resolve challenges related to assumptions
527 of spatial independence in species’ distributions and indirect effects of multiple species on associations.
528 Variations in the implementation of our approach can be used to test for association patterns that do not
529 follow the expectation from broad-scale climate, dispersal or other hypothesized processes. With network
530 metrics and analyses we can test for drivers of the patterns, and the approach can be potentially extended to
531 predictions of local community assembly. Combined with field experiments and other methods, our
532 framework can be a powerful tool to move beyond extant concepts in the analysis of co-occurrence data to
533 improve assembly predictions across scales and help merge community ecology and biogeography.

534

535 *Acknowledgements* – This study was conducted as a part of the BIEN Working Group (Principal
536 Investigators: Brian J. Enquist, Brad Boyle, Richard Condit, Steven Dolins, Robert K. Peet, and Barbara M.
537 Thiers) supported by the National Centre for Ecological Analysis and Synthesis, a center funded by the
538 National Science Foundation (NSF Grant EF-0553768), the University of California, Santa Barbara, and the
539 State of California. The BIEN Working Group was also supported by The iPlant Collaborative (NSF Grant
540 DBI-0735191). We thank all the contributors for the invaluable data provided to the BIEN
541 (<http://bien.nceas.ucsb.edu/bien/people/data-contributors/>). We are grateful to Irena Šímová for preparing the
542 trait data and to Brad Boyle for guidance on SQL queries. We thank D. Harris and two anonymous reviewers
543 for constructive comments that improved this manuscript, in particular D. Harris for pointing us towards the
544 use of partial correlation methods. We further acknowledge support to NMH by an EliteForsk Award, the
545 Aarhus University Research Foundation, and the Villum Foundation. B. Blonder and JCS acknowledge
546 financial support from Centre for Informatics Research on Complexity in Ecology (CIRCE), funded by the
547 Aarhus University Research Foundation under the AU Ideas, the European Research Council (ERC Starting
548 Grant #310886 ‘HISTFUNC’, and the Danish Council for Independent Research – Natural Sciences (grant
549 #12-125079). B. Blonder was also supported by a United States National Science Foundation graduate
550 research fellowship. CV was supported by a Marie Curie International Outgoing Fellowship within the 7th
551 European Community Framework Program (DiversiTraits project, no. 221060).

552 **References**

- 553 Afkhami, M. E. et al. 2014. Mutualist-mediated effects on species' range limits across large geographic
554 scales. - *Ecol. Lett.* 17: 1265–1273.
- 555 Araújo, M. B. and Rozenfeld, A. 2014. The geographic scaling of biotic interactions. - *Ecography* 37: 1–10.
- 556 Araújo, M. B. et al. 2011. Using species co-occurrence networks to assess the impacts of climate change. -
557 *Ecography* 34: 897–908.
- 558 Azaele, S. et al. 2010. Inferring plant ecosystem organization from species occurrences. - *J. Theor. Biol.* 262:
559 323–329.
- 560 Bascompte, J. 2010. Structure and dynamics of ecological networks. - *Science* 329: 765–766.
- 561 Baselga, A. 2010. Partitioning the turnover and nestedness components of beta diversity. - *Glob. Ecol.*
562 *Biogeogr.* 19: 134–143.
- 563 Baselga, A. 2012. The relationship between species replacement, dissimilarity derived from nestedness, and
564 nestedness. - *Glob. Ecol. Biogeogr.* 21: 1223–1232.
- 565 Benjamini, Y. and Hochberg, Y. 1995. Controlling the false discovery rate: A practical and powerful
566 approach to multiple testing. - *J. R. Stat. Soc. Ser. B* 57: 289–300.
- 567 Blick, R. and Burns, K. C. 2009. Network properties of arboreal plants: Are epiphytes, mistletoes and lianas
568 structured similarly? - *Perspect. Plant Ecol. Evol. Syst.* 11: 41–52.
- 569 Blois, J. L. et al. 2014. A framework for evaluating the influence of climate, dispersal limitation, and biotic
570 interactions using fossil pollen associations across the late Quaternary. - *Ecography* 37: 1095–1108.
- 571 Bolnick, D. I. et al. 2011. Why intraspecific trait variation matters in community ecology. - *Trends Ecol.*
572 *Evol.* 26: 183–192.

- 573 Borthagaray, A. I. et al. 2014. Inferring species roles in metacommunity structure from species co-
574 occurrence networks. - Proc. Biol. Sci. 281: 20141425.
- 575 Boulangeat, I. et al. 2012. Accounting for dispersal and biotic interactions to disentangle the drivers of
576 species distributions and their abundances. - Ecol. Lett. 15: 584–593.
- 577 Brown, J. H. et al. 1996. The geographic range: size, shape, boundaries, and internal structure. - Annu. Rev.
578 Ecol. Syst. 27: 597–623.
- 579 Brown, J. H. et al. 2004. Constraints on negative relationships. - In: Taper, M. L. and Lele, S. R. (eds), The
580 nature of scientific evidence: Statistical, philosophical, and empirical considerations. University of
581 Chicago Press, pp. 298–324.
- 582 Bruelheide, H. 2000. A new measure of fidelity and its application to defining species groups. - J. Veg. Sci.
583 11: 167–178.
- 584 Bullock, J. M. et al. 2000. Geographical separation of two *Ulex* species at three spatial scales: does
585 competition limit species' ranges? - Ecography 23: 257–271.
- 586 Calabrese, J. M. et al. 2014. Stacking species distribution models and adjusting bias by linking them to
587 macroecological models. - Glob. Ecol. Biogeogr. 23: 99–112.
- 588 Callaway, R. M. and Pennings, S. C. 2000. Facilitation may buffer competitive effects: Indirect and diffuse
589 interactions among salt marsh plants. - Am. Nat. 156: 416–424.
- 590 Callaway, R. M. et al. 2002a. Positive interactions among alpine plants increase with stress. - Nature 417:
591 844–848.
- 592 Callaway, R. M. et al. 2002b. Epiphyte host preferences and host traits: mechanisms for species-specific
593 interactions. - Oecologia 132: 221–230.
- 594 Chesson, P. 2000. Mechanisms of maintenance of species diversity. - Annu. Rev. Ecol. Syst. 31: 343–366.

595 Clauset, A. et al. 2004. Finding community structure in very large networks. - Phys. Rev. E 70: 066111.

596 Connor, E. F. and Simberloff, D. 1979. The assembly of species communities: chance or competition? -
597 Ecology 60: 1132–1140.

598 De Marco, P. et al. 2008. Spatial analysis improves species distribution modelling during range expansion. -
599 Biol. Lett. 4: 577–580.

600 Diamond, J. M. 1975. Assembly of species communities. - In: Cody, M. L. and Diamond, J. M. (eds),
601 Harvard University Press, pp. 342–444.

602 Dobra, A. et al. 2004. Sparse graphical models for exploring gene expression data. - J. Multivar. Anal. 90:
603 196–212.

604 Enquist, B. J. et al. 2009. The Botanical Information and Ecology Network (BIEN): Cyberinfrastructure for
605 an integrated botanical information network to investigate the ecological impacts of global climate change
606 on plant biodiversity. - iPlant Collaborative, www.iplantcollaborative.org

607 Erdős, P. and Rényi, A. 1959. On Random Graphs I. - Publ. Math. 6: 290 – 297.

608 Friedman, J. et al. 2008. Sparse inverse covariance estimation with the graphical lasso. - Biostatistics 9: 432–
609 441.

610 Gastauer, M. and Meira-Neto, J. A. A. 2013. Avoiding inaccuracies in tree calibration and phylogenetic
611 community analysis using Phylocom 4.2. - Ecol. Inform. 15: 85–90.

612 Goodman, L. A. and Kruskal, W. H. 1972. Measures of association for cross classifications, IV:
613 simplification of asymptotic variances. - J. Am. Stat. Assoc. 67: 415–421.

614 Gotelli, N. J. and Ulrich, W. 2010. The empirical Bayes approach as a tool to identify non-random species
615 associations. - Oecologia 162: 463–477.

616 Gotelli, N. J. et al. 2010. Macroecological signals of species interactions in the Danish avifauna. - Proc. Natl.
617 Acad. Sci. 107: 5030–5035.

618 Gray, A. et al. 2012. Forest Inventory and Analysis Database of the United States of America (FIA) (J
619 Dengler, J Oldeland, F Jansen, C M, E J, M Finckh, G F, G Lopez-Gonzalez, RK Peet, and JHJ
620 Schaminee, Eds.). - Biodivers. Ecol. 4: 225–231.

621 Guisan, A. and Zimmermann, N. E. 2000. Predictive habitat distribution models in ecology. - Ecol. Modell.
622 135: 147–186.

623 Guisan, A. and Rahbek, C. 2011. SESAM - a new framework integrating macroecological and species
624 distribution models for predicting spatio-temporal patterns of species assemblages. - J. Biogeogr. 38:
625 1433–1444.

626 Gutiérrez, E. E. et al. 2014. Can biotic interactions cause allopatry? Niche models, competition, and
627 distributions of South American mouse opossums. - Ecography 37: 741–753.

628 Harris, D. J. 2015. Estimating species interactions from observational data with Markov networks. - bioRxiv
629 [dx.doi.org/10.1101/018861](https://doi.org/10.1101/018861).

630 Hijmans, R. J. et al. 2005. Very high resolution interpolated climate surfaces for global land areas. - Int. J.
631 Climatol. 25: 1965–1978.

632 HilleRisLambers, J. et al. 2012. Rethinking community assembly through the lens of coexistence theory. -
633 Annu. Rev. Ecol. Evol. Syst. 43: 227–248.

634 Hoerl, A. E. and Kennard, R. W. 1970. Ridge regression: biased estimation for nonorthogonal problems. -
635 Technometrics 12: 55–67.

636 Janson, S. and Vegelius, J. 1981. Measures of ecological association. - Oecologia 49: 371–376.

637 Jongejans, E. et al. 2008. Dispersal, demography and spatial population models for conservation and control
638 management. - Perspect. Plant Ecol. Evol. Syst. 9: 153–170.

- 639 Kissling, W. D. et al. 2012. Towards novel approaches to modelling biotic interactions in multispecies
640 assemblages at large spatial extents. - *J. Biogeogr.* 39: 2163–2178.
- 641 Lasky, J. R. et al. 2014. Trait-mediated assembly processes predict successional changes in community
642 diversity of tropical forests. - *Proc. Natl. Acad. Sci.* 111: 5616–5621.
- 643 Lennon, J. J. 2000. Red-shifts and red herrings in geographical ecology. - *Ecography* 23: 101–113.
- 644 Leger, J.-B. et al. 2015. Clustering methods differ in their ability to detect patterns in ecological networks (T
645 Münkemüller, Ed.). - *Methods Ecol. Evol.* 6: 474–481.
- 646 Legendre, P. 1993. Spatial autocorrelation: trouble or new paradigm? - *Ecology* 74: 1659–1673.
- 647 Liaw, A. and Wiener, M. 2002. Classification and regression by randomForest. - *R News* 2: 18 – 22.
- 648 MacArthur, R. H. 1958. Population Ecology of Some Warblers of Northeastern Coniferous Forests. -
649 *Ecology* 39: 599.
- 650 Martorell, C. and Freckleton, R. P. 2014. Testing the roles of competition, facilitation and stochasticity on
651 community structure in a species-rich assemblage (R Brooker, Ed.). - *J. Ecol.* 102: 74–85.
- 652 Morales-Castilla, I. et al. 2015. Inferring biotic interactions from proxies. - *Trends Ecol. Evol.* 30: 347–356.
- 653 Newman, M. 2010. *Networks: an introduction.* - Oxford University Press.
- 654 Normand, S. et al. 2011. Postglacial migration supplements climate in determining plant species ranges in
655 Europe. - *Proc. R. Soc. B Biol. Sci.* 278: 3644–3653.
- 656 Olito, C. and Fox, J. W. 2015. Species traits and abundances predict metrics of plant-pollinator network
657 structure, but not pairwise interactions. - *Oikos* 124: 428–436.
- 658 Pastor, J. and Mladenoff, D. J. 1992. The southern boreal-northern hardwood forest border. - In: Shugart, H.
659 H. et al. (eds), *A systems analysis of the global boreal forest.* Cambridge University Press, pp. 216–240.

- 660 Phillips, S. J. and Dudík, M. 2008. Modeling of species distributions with Maxent: new extensions and a
661 comprehensive evaluation. - *Ecography* 31: 161–175.
- 662 Pollock, L. J. et al. 2014. Understanding co-occurrence by modelling species simultaneously with a Joint
663 Species Distribution Model (JSDM). - *Methods Ecol. Evol.* 5: 397–406.
- 664 Pottier, J. et al. 2013. The accuracy of plant assemblage prediction from species distribution models varies
665 along environmental gradients (R Field, Ed.). - *Glob. Ecol. Biogeogr.* 22: 52–63.
- 666 Schäfer, J. and Strimmer, K. 2005. A shrinkage approach to large-scale covariance matrix estimation and
667 implications for functional genomics. - *Stat. Appl. Genet. Mol. Biol.* 4: Article 32.
- 668 Schleuning, M. et al. 2012. Specialization of mutualistic interaction networks decreases toward tropical
669 latitudes. - *Curr. Biol.* 22: 1925–1931.
- 670 Soberón, J. 2007. Grinnellian and Eltonian niches and geographic distributions of species. - *Ecol. Lett.* 10:
671 1115–1123.
- 672 Somers, R. H. 1962. A new asymmetric measure of association for ordinal variables. - *Am. Sociol. Rev.* 27:
673 799–811.
- 674 Stewart, R. E. and Aldrich, J. W. 1951. Removal and Repopulation of Breeding Birds in a Spruce-Fir Forest
675 Community. - *Auk* 68: 471–482.
- 676 Svenning, J.-C. and Skov, F. 2007. Could the tree diversity pattern in Europe be generated by postglacial
677 dispersal limitation? - *Ecol. Lett.* 10: 453–460.
- 678 ter Steege, H. et al. 2013. Hyperdominance in the Amazonian tree flora. - *Science* 342: 1243092.
- 679 Thuiller, W. et al. 2013. A road map for integrating eco-evolutionary processes into biodiversity models. -
680 *Ecol. Lett.* 16 Suppl 1: 94–105.

- 681 Ulrich, W. et al. 2012. Null model tests for niche conservatism, phylogenetic assortment and habitat filtering.
682 - *Methods Ecol. Evol.* 3: 930–939.
- 683 Veech, J. A. 2013. A probabilistic model for analysing species co-occurrence (P Peres-Neto, Ed.). - *Glob.*
684 *Ecol. Biogeogr.* 22: 252–260.
- 685 Violle, C. et al. 2011. Phylogenetic limiting similarity and competitive exclusion. - *Ecol. Lett.* 14: 782–787.
- 686 Webb, C. O. et al. 2002. Phylogenies and community ecology. - *Annu. Rev. Ecol. Syst.* 33: 475–505.
- 687 Webb, C. O. et al. 2008. Phylocom: software for the analysis of phylogenetic community structure and trait
688 evolution. - *Bioinformatics* 24: 2098–2100.
- 689 Weiher, E. et al. 2011. Advances, challenges and a developing synthesis of ecological community assembly
690 theory. - *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* 366: 2403–2413.
- 691 Westoby, M. et al. 2002. Plant ecological strategies: Some leading dimensions of variation between species.
692 - *Annu. Rev. Ecol. Syst.* 33: 125–159.
- 693 Wisz, M. S. et al. 2013. The role of biotic interactions in shaping distributions and realised assemblages of
694 species: implications for species distribution modelling. - *Biol. Rev. Camb. Philos. Soc.* 88: 15–30.
- 695 Yuan, Y. et al. 2011. Directed partial correlation: inferring large-scale gene regulatory network through
696 induced topology disruptions. - *PLoS One* 6: e16835.

697 **Tables and Figures**

698 **Table 1.** Descriptive statistics for association networks constructed from the same dataset using different null
 699 models. We report an overall p-value for deviations from a Poisson distribution of edges, and for the subsets
 700 of the network for positive or negative associations, the mean \pm standard deviation of degree, as well as the
 701 number of non-singleton modules (i.e. with more than one member) detected in the network.

702

Regional model	P-value	Degree		Positive modules		Negative modules	
		Positive	Negative	Number	Mean size	Number	Mean size
Swap	2.44E-46	20.31 \pm 10.81	56.26 \pm 15.42	12	11.4	5	27.4
LOESS	0.230	0.09 \pm 0.31	0 \pm 0	5	2.2	0	NA
MaxEnt	0.565	0.82 \pm 0.89	0 \pm 0	24	3.2	0	NA

703

704

705 **Figure captions**

706 **Box 1.** Concepts and data flow underlying the framework. Each step illustrates how the network of non-
 707 random species association is derived from the observed co-occurrence matrix and the expectation based on
 708 a chosen null model. See detailed description in main text.

709

710 **Box 2.** Examples of possible analyses to explore species associations based on the network. Network-level
 711 analyses (left) give information on the overall structure of the network. Examples include 1) testing for
 712 overall-nonrandom number of links comparing the degree distribution to e.g. a binomial distribution, and 2)
 713 quantifying modularity, i.e. groups of species more likely to be associated. Node-level analyses (right) serve
 714 to 1) identify the roles of individual species (e.g. as aggregators or segregators by looking at the proportion
 715 of negative and positive links per species), 2) quantify the centrality of each species, or 3) test for
 716 correlations of node metrics such as network distance and hypothesized drivers such as phylogenetic or
 717 functional trait distance. Finally, a hybrid of network and node-level measures can be used for spatial

718 analyses (bottom) to test for trends in e.g. mean degree of communities across space and their correlation to
719 climate gradients.

720

721 **Figure 1.** Empirical association network for North American trees. The network was constructed using the
722 MaxEnt regional null model and a shrinkage inverse covariance estimator using 1,000 null replicates and a
723 false discovery rate of $\alpha=0.05$. Gray envelopes indicate distinct modules. Positive associations are shown as
724 blue lines; negative would be shown as red lines but none were found.

725

726 **Figure S1.** Association networks constructed from the same data, using the random-swap algorithm, the
727 LOESS, and the MaxEnt regional null models, respectively. Left panel, positive associations (blue); right
728 panel, negative associations (red). Gray envelopes indicate distinct modules.

729

730 **Figure S2.** Ratio of positive and negative associations across networks. The role of a species in the
731 network can be indicated by the degree of each node (i.e. number of associations) for the positive
732 subset of the network relative to the degree of each node for the negative subset of the network.
733 Species below the 1:1 line (purple) are ‘aggregators’ because they have proportionately more
734 positive associations, while those above the line are ‘segregators’. This designation shifts between
735 methods, as does the identity of the most highly associated species. The five species with highest
736 degree are labelled in each panel.

737

738 **Figure S3.** Trait and phylogenetic predictors of network structure. Points indicate pairwise
739 distances between all pairs of species. Relationships are shown for the random-swap network
740 (pink), LOESS network (orange), and MaxEnt network (purple).

741

742 **Figure S4.** Distribution of network statistics across space. Each point indicates the abundance-
743 weighted mean network degree for a different community.