

Food web reconstruction through phylogenetic transfer of low-rank network representation

Tanya Strydom^{1,2,‡} Salomé Bouskila^{1,‡} Francis Banville^{1,3,2} Ceres Barros⁴ Dominique Caron^{5,2}
Maxwell J Farrell⁶ Marie-Josée Fortin⁶ Victoria Hemming⁷ Benjamin Mercier^{3,2} Laura
J. Pollock^{5,2} Rogini Runghen⁸ Giulio V. Dalla Riva⁹ Timothée Poisot^{1,2}

¹ Département de Sciences Biologiques, Université de Montréal, Montréal, Canada ² Quebec Centre for Biodiversity Science, Montréal, Canada ³ Département de Biologie, Université de Sherbrooke, Sherbrooke, Canada ⁴ Department of Forest Resources Management, University of British Columbia, Vancouver, B.C., Canada ⁵ Department of Biology, McGill University, Montréal, Canada ⁶ Department of Ecology & Evolutionary Biology, University of Toronto, Toronto, Canada ⁷ Department of Forest and Conservation Sciences, University of British Columbia, Vancouver, Canada ⁸ Centre for Integrative Ecology, School of Biological Sciences, University of Canterbury, Canterbury, New Zealand ⁹ School of Mathematics and Statistics, University of Canterbury, Canterbury, New Zealand

[‡] These authors contributed equally to the work

Correspondance to:

Timothée Poisot — timothee.poisot@umontreal.ca

1. Despite their importance in many ecological processes, collecting data and information on ecological interactions is an exceedingly challenging task. For this reason, large parts of the world have a data deficit when it comes to species interactions, and how the resulting networks are structured. As data collection alone is unlikely to be sufficient, community ecologists must adopt predictive methods.
2. We present a methodological framework that uses graph embedding and transfer learning to assemble a predicted list of trophic interactions of a species pool for which their interactions are unknown. Specifically, we ‘learn’ the information from a known interaction network by inferring the latent traits of species and infer the latent traits of a species pool for which we have no *a priori* interaction data based on their phylogenetic relatedness to species from the known network. The latent traits can then be used to predict interactions and construct an interaction network.
3. Here we assembled a metaweb for Canadian mammals derived from interactions in the European food web, despite only 4% of common species being shared between the two locations. The results of the predictive model are compared against databases of recorded pairwise interactions, showing that we correctly recover 91% of known interactions.
4. The framework itself is robust even when the known network is incomplete or contains spurious interactions making it an ideal candidate as a tool for filling gaps when it comes to species interactions. We provide guidance on how this framework can be adapted by substituting some approaches or predictors in order to make it more generally applicable.

1 Introduction

2 There are two core challenges we are faced with in furthering our understanding of ecological networks
3 across space, particularly at macro-ecologically relevant scales (e.g. Trøjelsgaard & Olesen, 2016). First,
4 ecological networks within a location are difficult to sample properly (Jordano, 2016a, 2016b), resulting in
5 a widespread “Eltonian shortfall” (Hortal et al., 2015), *i.e.* a lack of knowledge about inter- and intra-
6 specific relationships. This first challenge has been, in large part, addressed by the recent emergence of a
7 suite of methods aiming to predict interactions within *existing* networks, many of which are reviewed in
8 Strydom, Catchen, et al. (2021). Second, recent analyses based on collected data (Poisot, Bergeron, et al.,
9 2021) or metadata (Cameron et al., 2019) highlight that ecological networks are currently studied in a
10 biased subset of space and bioclimates, which impedes our ability to generalize any local understanding of
11 network structure. Meaning that, although the framework to address incompleteness *within* networks
12 exists, there would still be regions for which, due to a *lack* of local interaction data, we are unable to infer
13 potential species interactions. Having a general solution for inferring *potential* interactions (despite the
14 unavailability of interaction data) could be the catalyst for significant breakthroughs in our ability to start
15 thinking about species interaction networks over large spatial scales. In a recent overview of the field of
16 ecological network prediction, Strydom, Catchen, et al. (2021) identified two challenges of interest to the
17 prediction of interactions at large scales. First, there is a relative scarcity of relevant data in most places
18 globally – paradoxically, this restricts our ability to infer interactions to locations where inference is
19 perhaps the least required; second, accurate predictions often demand accurate predictors, and the lack of
20 methods that can leverage small amount of data is a serious impediment to our predictive ability globally.

21 Here, we present a general method to infer potential trophic interactions, relying on the transfer learning
22 of network representations, specifically by using similarities of species in a biologically/ecologically
23 relevant proxy space (e.g. shared morphology or ancestry). Transfer learning is a machine learning
24 methodology that uses the knowledge gained from solving one problem and applying it to a related
25 (destination) problem (Pan & Yang, 2010; Torrey & Shavlik, 2010). In this instance, we solve the problem
26 of predicting trophic interactions between species, based on knowledge extracted from another species
27 pool for which interactions are known by using phylogenetic structure as a medium for transfer. There is a
28 plurality of measures of species similarities that can be used for inferring *potential* species interactions *i.e.*
29 metaweb reconstruction (see *e.g.* Morales-Castilla et al., 2015); however, phylogenetic proximity has

several desirable properties when working at large scales. Gerhold et al. (2015) made the point that phylogenetic signal captures diversification of characters (large macro-evolutionary process), but not necessarily community assembly (fine ecological process); Dormann et al. (2010) previously found very similar conclusions. Interactions tend to reflect a phylogenetic signal because they have a conserved pattern of evolutionary convergence that encompasses a wide range of ecological and evolutionary mechanisms (Cavender-Bares et al., 2009; Mouquet et al., 2012), and - most importantly - retain this signal even if it is obscured at the community scale due to e.g. local conditions (Hutchinson et al., 2017; Poisot & Stouffer, 2018). Finally, species interactions at macro-ecological scales seem to respond mostly to macro-evolutionary processes (Price, 2003); which is evidenced by the presence of conserved backbones in food webs (Dalla Riva & Stouffer, 2016; Mora et al., 2018), strong evolutionary signature on prey choice (Stouffer et al., 2012), and strong phylogenetic signature in food web intervalty (Eklöf & Stouffer, 2016). Phylogenetic reconstruction has also previously been used within the context of ecological networks, namely understanding ancestral plant-insect interactions (Braga et al., 2021). Taken together, these considerations suggest that phylogenies can reliably be used to transfer knowledge on species interactions.

[Figure 1 about here.]

In fig. 1, where we provide a methodological overview based on learning the embedding of a metaweb of trophic interactions for European mammals (known interactions; Maiorano et al., 2020a, 2020b) and, based on phylogenetic relationships between mammals globally (*i.e.*, phylogenetic tree Upham et al., 2019), infer a metaweb for the Canadian mammalian species pool (using only a species list *i.e.* interactions are ‘unknown’ in this instance). Following the definition of Dunne (2006), a metaweb is a network analogue to the regional species pool; specifically, it is an inventory of all *potential* interactions between species from a spatially delimited area (and so captures the γ diversity of interactions). The metaweb is, therefore, *not* a prediction of the food web at a specific locale within the spatial area it covers, and will have a different structure (notably by having a larger connectance; see *e.g.* Wood et al., 2015). These local food webs (which captures the α diversity of interactions) are a subset of the metaweb’s species and interactions, and have been called “metaweb realizations” (Poisot et al., 2015). Differences between local food web and their metaweb are due to chance, species abundance and co-occurrence, local environmental conditions, and local distribution of functional traits, among others.

Because the metaweb represents the joint effect of functional, phylogenetic, and macroecological

59 processes (Morales-Castilla et al., 2015), it holds valuable ecological information. Specifically, it is the
60 “upper bounds” on what the composition of the local networks can be (see e.g. McLeod et al., 2021). These
61 local networks, in turn, can be reconstructed given appropriate knowledge of local species composition,
62 providing information on structure of food webs at finer spatial scales. This has been done for example for
63 tree-galler-parasitoid systems (Gravel et al., 2018), fish trophic interactions (Albouy et al., 2019), tetrapod
64 trophic interactions (O’Connor et al., 2020), and crop-pest networks (Grünig et al., 2020). Whereas the
65 original metaweb definition, and indeed most past uses of metawebs, was based on the presence/absence
66 of interactions, we focus on *probabilistic* metawebs where interactions are represented as the chance of
67 success of a Bernoulli trial (see e.g. Poisot et al., 2016); therefore, not only does our method recommend
68 interactions that may exist, it gives each interaction a score, allowing us to properly weigh them.

69 Our case study shows that phylogenetic transfer learning is an effective approach to the generation of
70 probabilistic metawebs. This showcases that although the components (species) that make up the
71 Canadian and European communities may be *minimally* shared (the overall species overlap is less than
72 4%), if the medium (proxy space) selected in the transfer step is biologically plausible, we can still
73 effectively learn from the known network and make biologically relevant predictions of interactions.

74 Indeed, as we detail in the results, when validated against known but fractional data of trophic
75 interactions between Canadian mammals, our model achieves a predictive accuracy of approximately 91%.
76 It should be reiterated that the framework presented in fig. 1 is amenable to changes e.g. it is possible to
77 use distinct trees if working with distinct clades (such as pollination networks) or, alternatively, the
78 measure of similarity may not be phylogeny, and can be replaced by information on foraging (Beckerman
79 et al., 2006), cell-level mechanisms (Boeckaerts et al., 2021), or a combination of traits and phylogenetic
80 structure (Stock, 2021). Most importantly, although we focus on a trophic system, it is an established fact
81 that different (non-trophic) interactions do themselves interact with and influence the outcome of trophic
82 interactions (see e.g. Kawatsu et al., 2021; Kéfi et al., 2012). Future development of metaweb inference
83 techniques should cover the prediction of multiple interaction types.

84 Data used for the case study

85 We use data from the European metaweb assembled by Maiorano et al. (2020a). This was assembled using
86 data extracted from scientific literature (including published papers, books, and grey literature) from the

87 last 50 years and includes all terrestrial tetrapods (mammals, breeding birds, reptiles and amphibians)
88 occurring on the European sub-continent (and Turkey) - with the caveat that only species introduced in
89 historical times and currently naturalized being included. The European metaweb was filtered using the
90 Global Biodiversity Information Facility (GBIF) taxonomic backbone (GBIF Secretariat, 2021) so as to
91 contain only terrestrial and semi-aquatic mammals. As all species had valid matches to the GBIF
92 taxonomy it was used as the backbone for the remaining reconciliation steps namely, the mammalian
93 consensus supertree by Upham et al. (2019) (which is used for the knowledge transfer step) and for the
94 Canadian species list—which was extracted from the International Union for Conservation of Nature
95 (IUCN) checklist, and corresponds to the same selection criteria that was applied by Maiorano et al.
96 (2020a) in the European metaweb.

97 After taxonomic cleaning and reconciliation the mammalian European metaweb has 260 species, and the
98 Canadian species pool 163; of these, 17 (about 4% of the total) are shared, and 89 species from Canada
99 (54%) had at least one congeneric species in Europe. The similarity for both species pools predictably
100 increases with higher taxonomic order, with 19% of shared genera, 47% of shared families, and 75% of
101 shared orders; for the last point, Canada and Europe each had a single unique order (*Didelphimorphia* for
102 Canada, *Erinaceomorpha* for Europe).

103 **Method description**

104 The core point of our method is the transfer of knowledge of a known ecological network, in order to
105 predict interactions between species from another location for which the network is unknown (or partially
106 known). In fig. 1 the grey text boxes give a high-level overview of the approach; in the example around
107 which this manuscript is built The method we develop is, ecologically speaking, a “black box,” i.e. an
108 algorithm that can be understood mathematically, but whose component parts are not always directly tied
109 to ecological processes. There is a growing realization in machine learning that (unintentional) black box
110 algorithms are not necessarily a bad thing (Holm, 2019), as long as their constituent parts can be
111 examined (which is the case with our method). But more importantly, data hold more information than
112 we might think; as such, even algorithms that are disconnected from a model can make correct guesses
113 most of the time (Halevy et al., 2009); in fact, in an instance of ecological forecasting of spatio-temporal
114 systems, model-free approaches (i.e. drawing all of their information from the data) outperformed

115 model-informed ones (Perretti et al., 2013).

116 **Implementation and code availability**

117 The entire pipeline is implemented in *Julia* 1.6 (Bezanson et al., 2017) and is available under the
118 permissive MIT License at <https://osf.io/2zwqm/>. The taxonomic cleanup steps are done using GBIF.jl
119 (Dansereau & Poisot, 2021). The network embedding and analysis is done using EcologicalNetworks.jl
120 (Banville et al., 2021; Poisot et al., 2019). The phylogenetic simulations are done using PhyloNetworks.jl
121 (Solís-Lemus et al., 2017) and Phylo.jl (Reeve et al., 2016). A complete Project.toml file specifying the
122 full tree of dependencies is available alongside the code. This material also includes a fully annotated copy
123 of the entire code required to run this project (describing both the intent of the code and discussing some
124 technical implementation details), a vignette for every step of the process, and a series of Jupyter
125 notebooks with the text and code. The pipeline can be executed on a laptop in a matter of minutes, and
126 therefore does not require extensive computational power.

127 **Step 1: Learning the origin network representation**

128 The first step in transfer learning is to learn the structure of the original dataset. In order to do so, we rely
129 on an approach inspired from representational learning, where we learn a *representation* of the metaweb
130 (in the form of the latent subspaces), rather than a list of interactions (species *a* eats *b*). This approach is
131 conceptually different from other metaweb-scale predictions (e.g. Albouy et al., 2019), in that the metaweb
132 representation is easily transferable. Specifically, we use Random Dot Product Graph model (hereafter
133 RDPG; S. J. Young & Scheinerman, 2007) to create a number of latent variables that can be combined into
134 an approximation of the network adjacency matrix. RDPG is known to capture the evolutionary backbone
135 of food webs (Dalla Riva & Stouffer, 2016), resulting in strong phylogenetic signal in RDPG results; in
136 other words, the latent variables of an RDPG can be mapped onto a phylogenetic tree, and
137 phylogenetically similar predators should share phylogenetically similar preys. In addition, recent
138 advances show that the latent variables produced this way can be used to predict *de novo* interactions.
139 Interestingly, the latent variables do not need to be produced by decomposing the network itself; in a
140 recent contribution, Runghen et al. (2021) showed that deep artificial neural networks are able to
141 reconstruct the left and right subspaces of an RDPG, in order to predict human movement networks from

142 individual/location metadata and opens up the possibility of using additional metadata as predictors.

143 The latent variables are created by performing a truncated Singular Value Decomposition (t-SVD; Halko et
144 al., 2011) on the adjacency matrix. SVD is an appropriate embedding of ecological networks, which has
145 recently been shown to both capture their complex, emerging properties (Strydom, Dalla Riva, et al., 2021)
146 and to allow highly accurate prediction of the interactions within a single network (Poisot, Ouellet, et al.,
147 2021). Under SVD, an adjacency matrix \mathbf{A} (where $\mathbf{A}_{m,n} \in \mathbb{B}$ where 1 indicates predation and 0 an absence
148 thereof) is decomposed into three components resulting in $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}'$. Here, Σ is a $m \times n$ diagonal matrix
149 and contains only singular (σ) values along its diagonal, \mathbf{U} is a $m \times m$ unitary matrix, and \mathbf{V}' a $n \times n$
150 unitary matrix. Truncating the SVD removes additional noise in the dataset by omitting non-zero and/or
151 smaller σ values from Σ using the rank of the matrix. Under a t-SVD $\mathbf{A}_{m,n}$ is decomposed so that Σ is a
152 square $r \times r$ diagonal matrix (with $1 \leq r \leq r_{full}$ where r_{full} is the full rank of \mathbf{A} and r the rank at which we
153 truncate the matrix) containing only non-zero σ values. Additionally, \mathbf{U} is now a $m \times r$ semi-unitary
154 matrix and \mathbf{V}' a $n \times r$ semi-unitary matrix.

155 The specific rank at which the SVD ought to be truncated is a difficult question. The purpose of SVD is to
156 remove the noise (expressed at high dimensions) and to focus on the signal, (expressed at low dimensions).
157 In datasets with a clear signal/noise demarcation, a scree plot of Σ can show a sharp drop at the rank where
158 noise starts (Zhu & Ghodsi, 2006). Because the European metaweb is almost entirely known, the amount
159 of noise (uncertainty) is low; this is reflected in fig. 2 (left), where the scree plot shows no important drop,
160 and in fig. 2 (right) where the proportion of variance explained increases smoothly at higher dimensions.
161 For this reason, we default back to a threshold that explains 60% of the variance in the underlying data,
162 corresponding to 12 dimensions - *i.e.* a tradeoff between accuracy and a reduced number of features.

163 An RDPG estimates the probability of observing interactions between nodes (species) as a function of the
164 nodes' latent variables, and is a way to turn an SVD (which decompose one matrix into three) into two
165 matrices that can be multiplied to provide an approximation of the network. The latent variables used for
166 the RDPG, called the left and right subspaces, are defined as $\mathcal{L} = \mathbf{U}\sqrt{\Sigma}$, and $\mathcal{R} = \sqrt{\Sigma}\mathbf{V}'$ – using the full
167 rank of \mathbf{A} , $\mathcal{L}\mathcal{R} = \mathbf{A}$, and using any smaller rank results in $\mathcal{L}\mathcal{R} \approx \mathbf{A}$. Using a rank of 1 for the t-SVD
168 provides a first-order approximation of the network. One advantage of using an RDPG rather than an SVD
169 is that the number of components to estimate decreases; notably, one does not have to estimate the
170 singular values of the SVD. Furthermore, the two subspaces can be directly multiplied to yield a network.

172 Because RDPG relies on matrix multiplication, the higher dimensions essentially serve to make specific
173 interactions converge towards 0 or 1; therefore, for reasonably low ranks, there is no guarantee that the
174 values in the reconstructed network will be within the unit range. In order to determine what constitutes
175 an appropriate threshold for probability, we performed the RDPG approach on the European metaweb,
176 and evaluated the probability threshold by treating this as a binary classification problem, specifically
177 assuming that both 0 and 1 in the European metaweb are all true. Given the methodological details given
178 in Maiorano et al. (2020a) and O'Connor et al. (2020), this seems like a reasonable assumption, although
179 one that does not hold for all metawebs. We used the thresholding approach presented in Poisot, Ouellet,
180 et al. (2021), and picked a cutoff that maximized Youden's J statistic (a measure of the informedness
181 (trust) of predictions; Youden (1950)); the resulting cutoff was 0.22, and gave an accuracy above 0.99. In
182 Supp. Mat. 1, we provide several lines of evidence that using the entire network to estimate the threshold
183 does not lead to overfitting; that using a subset of species would yield the same threshold; that decreasing
184 the quality of the original data by adding or removing interactions would minimally affect the predictive
185 accuracy of RDPG applied to the European metaweb; and that the networks reconstructed from artificially
186 modified data are reconstructed with the correct ecological properties.

187 The left and right subspaces for the European metaweb, accompanied by the threshold for prediction,
188 represent the knowledge we seek to transfer. In the next section, we explain how we rely on phylogenetic
189 similarity to do so.

190 **Steps 2 and 3: Transfer learning through phylogenetic relatedness**

191 In order to transfer the knowledge from the European metaweb to the Canadian species pool, we
192 performed ancestral character estimation using a Brownian motion model, which is a conservative
193 approach in the absence of strong hypotheses about the nature of phylogenetic signal in the network
194 decomposition (Litsios & Salamin, 2012). This uses the estimated feature vectors for the European
195 mammals to create a state reconstruction for all species (conceptually something akin to a trait-based
196 mammalian phylogeny using latent generality and vulnerability traits) and allows us to impute the
197 missing (latent) trait data for the Canadian species that are not already in the European network; as we are
198 focused on predicting contemporary interactions, we only retained the values for the tips of the tree. We

assumed that all traits (*i.e.* the feature vectors for the left and right subspaces) were independent, which is a reasonable assumption as every trait/dimension added to the t-SVD has an *additive* effect to the one before it. Note that the Upham et al. (2019) tree itself has some uncertainty associated to inner nodes of the phylogeny. In this case study, we have decided to not propagate this uncertainty, as it would complexify the process. The Brownian motion algorithm returns the *average* value of the trait, and its upper and lower bounds. Because we do not estimate other parameters of the traits' distributions, we considered that every species trait is represented as a uniform distribution between these bounds. The choice of the uniform distribution was made because the algorithm returns a minimum and maximum point estimate for the value, and given this information, the uniform distribution is the one with maximum entropy. Had all mean parameters estimates been positive, the exponential distribution would have been an alternative, but this is not the case for the subspaces of an RDPG. In order to examine the consequences of the choice of distribution, we estimated the variance per latent variable per node to use a Normal distribution; as we show in Supp. Mat. 2, this decision results in dramatically over-estimating the number and probability of interactions, and therefore we keep the discussions in the main text to the uniform case. The inferred left and right subspaces for the Canadian species pool ($\hat{\mathcal{L}}$ and $\hat{\mathcal{R}}$) have entries that are distributions, representing the range of values for a given species at a given dimension.

These objects represent the transferred knowledge, which we can use for prediction of the Canadian metaweb.

Step 4: Probabilistic prediction of the destination network

The phylogenetic reconstruction of $\hat{\mathcal{L}}$ and $\hat{\mathcal{R}}$ has an associated uncertainty, represented by the breadth of the uniform distribution associated to each of their entries. Therefore, we can use this information to assemble a *probabilistic* metaweb in the sense of Poisot et al. (2016), *i.e.* in which every interaction is represented as a single, independent, Bernoulli event of probability p .

[Figure 3 about here.]

Specifically, we have adopted the following approach. For every entry in $\hat{\mathcal{L}}$ and $\hat{\mathcal{R}}$, we draw a value from its distribution. This results in one instance of the possible left ($\hat{\ell}$) and right (\hat{r}) subspaces for the Canadian metaweb. These can be multiplied, to produce one matrix of real values. Because the entries in

226 $\hat{\ell}$ and \hat{r} are in the same space where \mathcal{L} and \mathcal{R} were originally predicted, it follows that the threshold ρ
227 estimated for the European metaweb also applies. We use this information to produce one random
228 Canadian metaweb, $N = \hat{\mathcal{L}}\hat{\mathcal{R}}' \geq \rho$. As we can see in (fig. 3), the European and Canadian metawebs are
229 structurally similar (as would be expected given the biogeographic similarities) and the two (left and right)
230 subspaces are distinct *i.e.* capturing predation (generality) and prey (vulnerability) latent traits.

231 Because the intervals around some trait values can be broad (in fact, probably broader than what they
232 would actually be, see *e.g.* Garland et al., 1999), we repeat the above process 2×10^5 times, which results in
233 a probabilistic metaweb P , where the probability of an interaction (here conveying our degree of trust that
234 it exists given the inferred trait distributions) is given by the number of times where it appears across all
235 random draws N , divided by the number of samples. An interaction with $P_{i,j} = 1$ means that these two
236 species were predicted to interact in all 2×10^5 random draws.

237 It must be noted that despite bringing in a large amount of information from the European species pool
238 and interactions, the Canadian metaweb has distinct structural properties. Following an approach similar
239 to Vermaat et al. (2009), we show in Supp. Mat. 3 that not only can we observe differences in a
240 multivariate space between the European and Canadian metaweb, we can also observe differences in the
241 same space between random subgraphs from these networks. These results line up with the studies
242 spatializing metawebs that have been discussed in the introduction: changes in the species pool are
243 driving local structural changes in the networks.

244 **Data cleanup, discovery, validation, and thresholding**

245 Once the probabilistic metaweb for Canada has been produced, we followed a number of data inflation
246 steps to finalize it. This step is external to the actual transfer learning framework but rather serves as a
247 way to augment and validate the predicted metaweb.

248 [Figure 4 about here.]

249 First, we extracted the subgraph corresponding to the 17 species shared between the European and
250 Canadian pools and replaced these interactions with a probability of 0 (non-interaction) or 1 (interaction),
251 according to their value in the European metaweb. This represents a minute modification of the inferred

252 network (about 0.8% of all species pairs from the Canadian web), but ensures that we are directly re-using
253 knowledge from Europe.

254 Second, we looked for all species in the Canadian pool known to the Global Biotic Interactions (GloBI)
255 database (Poelen et al., 2014), and extracted their known interactions. Because GloBI aggregates observed
256 interactions, it is not a *networks* data source, and therefore the only information we can reliably extract
257 from it is that a species pair *was reported to interact at least once*. This last statement should yet be taken
258 with caution, as some sources in GloBI (e.g. Thessen & Parr, 2014) are produced through text analysis,
259 and therefore may not document direct evidence of the interaction. Nevertheless, should the predictive
260 model work, we would expect that a majority of interactions known to GloBI would also be predicted. We
261 retrieved 366 interactions between mammals from the Canadian species pool from GloBI, 33 of which
262 were not predicted by the model; this results in a success rate of 91%. After performing this check, we set
263 the probability of all interactions known to GloBI to 1.

264 Finally, we downloaded the data from Strong & Leroux (2014), who mined various literature sources to
265 identify trophic interactions in Newfoundland. This dataset documented 25 interactions between
266 mammals, only two of which were not part of our (Canada-level) predictions, resulting in a success rate of
267 92%. These two interactions were added to our predicted metaweb with a probability of 1. A table listing
268 all interactions in the predicted Canadian metaweb can be found in the supplementary material.

269 [Figure 5 about here.]

270 Because the confidence intervals on the inferred trait space are probably over-estimates, we decided to
271 apply a thresholding step to the interactions after the data inflation (see fig. 5 showing the effect of varying
272 the cutoff on $P(i \rightarrow j)$). Cirtwill & Hämäck (2021) proposed a number of strategies to threshold
273 probabilistic networks. Their methods assume the underlying data to be tag-based sequencing, which
274 represents interactions as co-occurrences of predator and prey within the same tags; this is conceptually
275 identical to our Bernoulli-trial based reconstruction of a probabilistic network. We performed a full
276 analysis of the effect of various cutoffs, and as they either resulted in removing too few interactions, or
277 removing enough interactions that species started to be disconnected from the network, we set this
278 threshold for a probability equivalent to 0 to the largest possible value that still allowed all species to have
279 at least one interaction with a non-zero probability. The need for this slight deviation from the Cirtwill &
280 Hämäck (2021) method highlights the need for additional development on network thresholding.

281 **Results and discussion of the case study**

282 [Figure 6 about here.]

283 The t-SVD embedding is able to learn relevant ecological features for the network. fig. 6 shows that the
284 first rank correlates linearly with generality and vulnerability (Schoener, 1989), *i.e.* the number of preys
285 and predators for each species. Importantly, this implies that a rank 1 approximation represents the
286 configuration model for the metaweb, *i.e.* a set of random networks generated from a given degree
287 sequence (Park & Newman, 2004). Accounting for the probabilistic nature of the degrees, the rank 1
288 approximation also represents the *soft* configuration model (van der Hoorn et al., 2018). Both models are
289 maximum entropy graph models (Garlaschelli et al., 2018), with sharp (all network realizations satisfy the
290 specified degree sequence) and soft (network realizations satisfy the degree sequence on average) local
291 constraints, respectively. The (soft) configuration model is an unbiased random graph model widely used
292 by ecologists in the context of null hypothesis significance testing of network structure (*e.g.* Bascompte et
293 al., 2003) and can provide informative priors for Bayesian inference of network structure (*e.g.* J.-G. Young
294 et al., 2021). It is noteworthy that for this metaweb, the relevant information was extracted at the first
295 rank. Because the first rank corresponds to the leading singular value of the system, the results of fig. 6
296 have a straightforward interpretation: degree-based processes are the most important in structuring the
297 mammalian food web.

298 **Discussion**

299 One important aspect in which Europe and Canada differ (despite their comparable bioclimatic
300 conditions) is the degree of the legacy of human impacts, which have been much longer in Europe.
301 Nenzén et al. (2014) showed that even at small scales (the Iberian peninsula), mammal food webs retain
302 the signal of both past climate change and human activity, even when this human activity was orders of
303 magnitude less important than it is now. Similarly, Yeakel et al. (2014) showed that changes in human
304 occupation over several centuries can lead to food web collapse. Megafauna in particular seems to be very
305 sensitive to human arrival (Pires et al., 2015). In short, there is well-substantiated support for the idea that
306 human footprint affects more than the risk of species extinction (Marco et al., 2018), and can lead to
307 changes in interaction structure. Yet, owing to the inherent plasticity of interactions, there have been

308 documented instances of food webs undergoing rapid collapse/recovery cycles over short periods of time
309 (Pedersen et al., 2017). The embedding of a network, in a sense, embeds its macro-evolutionary history,
310 especially as RDPG captures ecological signal (Dalla Riva & Stouffer, 2016); at this point, it is important to
311 recall that a metaweb is intended as a catalogue of all potential interactions, which should then be filtered
312 (Morales-Castilla et al., 2015). In practice (and in this instance) the reconstructed metaweb will predict
313 interactions that are plausible based on the species' evolutionary history, however some interactions
314 would/would not be realized due to human impact.

315 Dallas et al. (2017) suggested that most links in ecological networks may be cryptic, *i.e.* uncommon or
316 otherwise hard to observe. This argument essentially echoes Jordano (2016b): the sampling of ecological
317 interactions is difficult because it requires first the joint observation of two species, and then the
318 observation of their interaction. In addition, it is generally expected that weak or rare links would be more
319 common in networks (Csermely, 2004), compared to strong, persistent links; this is notably the case in
320 food chains, wherein many weaker links are key to the stability of a system (Neutel et al., 2002). In the
321 light of these observations, the results in fig. 4 are not particularly surprising: we expect to see a surge in
322 these low-probability interactions under a model that has a good predictive accuracy. Because the
323 predictions we generate are by design probabilistic, then one can weigh these rare links appropriately. In a
324 sense, that most ecological interactions are elusive can call for a slightly different approach to sampling:
325 once the common interactions are documented, the effort required in documenting each rare interaction
326 may increase exponentially. Recent proposals suggest that machine learning algorithms, in these
327 situations, can act as data generators (Hoffmann et al., 2019): in this perspective, high quality
328 observational data can be supplemented with synthetic data coming from predictive models, which
329 increases the volume of information available for inference. Indeed, Strydom, Catchen, et al. (2021)
330 suggested that knowing the metaweb may render the prediction of local networks easier, because it fixes
331 an “upper bound” on which interactions can exist; indeed, with a probabilistic metaweb, we can consider
332 that the metaweb represents an aggregation of informative priors on the interactions.

333 Related to the last point, Cirtwill et al. (2019) showed that network inference techniques based on
334 Bayesian approaches would perform far better in the presence of an interaction-level informative prior;
335 the desirable properties of such a prior would be that it is expressed as a probability, preferably
336 representing a Bernoulli event, the value of which would be representative of relevant biological processes
337 (probability of predation in this case). We argue that the probability returned at the very last step of our

338 framework may serve as this informative prior; indeed, the output of our analysis can be used in
339 subsequent steps, also possibly involving expert elicitation to validate some of the most strongly
340 recommended interactions. One important *caveat* to keep in mind when working with interaction
341 inference is that interactions can never really be true negatives (in the current state of our methodological
342 framework and data collection limitations); this renders the task of validating a model through the usual
343 application of binary classification statistics very difficult (although see Strydom, Catchen, et al., 2021 for a
344 discussion of alternative suggestions). The other way through which our framework can be improved is by
345 substituting the predictors that are used for transfer. For example, in the presence of information on
346 species traits that are known to be predictive of species interactions, one might want to rely on functional
347 rather than phylogenetic distances – in food webs, body size (and allometrically related variables) has
348 been established as such a variable (Brose et al., 2006); the identification of relevant functional traits is
349 facilitated by recent methodological developments (Rosado et al., 2013). It should be noted that Xing &
350 Fayle (2021) highlight phylogenetic relatedness as one of the core components of network comparison at
351 the global scale. In this case study, we have embedded the original metaweb using t-SVD, because it lends
352 itself to an RDPG reconstruction, which is known to capture the consequences of evolutionary processes
353 (Dalla Riva & Stouffer, 2016); this being said, there are other ways to embed graphs (Arsov & Mirceva,
354 2019; Cai et al., 2017; Cao et al., 2019), which can be used as alternatives.

355 As Herbert (1965) rightfully pointed out, “[y]ou can’t draw neat lines around planet-wide problems”; in
356 this regard, our approach (and indeed, any inference of a metaweb at large scales) must contend with
357 several interesting and interwoven families of problems. The first is the limit of the metaweb to embed
358 and transfer. If the initial metaweb is too narrow in scope, notably from a taxonomic point of view, the
359 chances of finding another area with enough related species to make a reliable inference decreases; this
360 would likely be indicated by large confidence intervals during ancestral character estimation, but the lack
361 of well documented metawebs is currently preventing the development of more concrete guidelines. The
362 question of phylogenetic relatedness and dispersal is notably true if the metaweb is assembled in an area
363 with mostly endemic species, and as with every predictive algorithm, there is room for the application of
364 our best ecological judgement. Conversely, the metaweb should be reliably filled, which assumes that the
365 S^2 interactions in a pool of S species have been examined, either through literature surveys or expert
366 elicitation. Supp. Mat. 1 provides some guidance as to the type of sampling effort that should be
367 prioritized. While RDPG was able to maintain very high predictive power when interactions were missing,

368 the addition of false positive interactions was immediately detected; this suggests that it may be
369 appropriate to err on the side of “too many” interactions when constructing the initial metaweb to be
370 transferred. The second series of problems are related to determining which area should be used to infer
371 the new metaweb in, as this determines the species pool that must be used. In our application, we focused
372 on the mammals of Canada. The upside of this approach is that information at the country level is likely
373 to be required by policy makers and stakeholders for their biodiversity assessment, as each country tends
374 to set goals at the national level (Buxton et al., 2021) for which quantitative instruments are designed
375 (Turak et al., 2017), with specific strategies often enacted at smaller scales (Ray et al., 2021). And yet, we
376 do not really have a satisfying answer to the question of “where does a food web stop?”; the current most
377 satisfying solutions involve examining the spatial consistency of network area relationships (Fortin et al.,
378 2021; see e.g. Galiana et al., 2018, 2019, 2021), which is of course impossible in the absence of enough
379 information about the network itself. This suggests that an *a posteriori* refinement of the results may be
380 required, based on a downscaling of the metaweb. The final family of problems relates less to the
381 availability of data or quantitative tools, and more to the praxis of spatial ecology. Operating under the
382 context of national divisions, in large parts of the world, reflects nothing more than the legacy of settler
383 colonialism. Indeed, the use of ecological data is not an apolitical act (Nost & Goldstein, 2021), as data
384 infrastructures tend to be designed to answer questions within national boundaries, and their use both
385 draws upon and reinforces territorial statecraft; as per Machen & Nost (2021), this is particularly true
386 when the output of “algorithmic thinking” (e.g. relying on machine learning to generate knowledge) can
387 be re-used for governance (e.g. enacting conservation decisions at the national scale). We therefore
388 recognize that methods such as we propose operate under the framework that contributed to the ongoing
389 biodiversity crisis (Adam, 2014), reinforced environmental injustice (Choudry, 2013; Domínguez &
390 Luoma, 2020), and on Turtle Island especially, should be replaced by Indigenous principles of land
391 management (Eichhorn et al., 2019; No’kmag et al., 2021). As we see AI/ML being increasingly mobilized
392 to generate knowledge that is lacking for conservation decisions (e.g. Lamba et al., 2019; Mosebo
393 Fernandes et al., 2020), our discussion of these tools need to go beyond the technical, and into the
394 governance consequences they can have.

395 **Acknowledgements:** We acknowledge that this study was conducted on land within the traditional
396 unceded territory of the Saint Lawrence Iroquoian, Anishinabewaki, Mohawk, Huron-Wendat, and
397 Omàmiwininiwak nations. TP, TS, DC, and LP received funding from the Canadian Institute for Ecology &

398 Evolution. FB is funded by the Institute for Data Valorization (IVADO). TS, SB, and TP are funded by a
399 donation from the Courtois Foundation. CB was awarded a Mitacs Elevate Fellowship no. IT12391, in
400 partnership with fRI Research, and also acknowledges funding from Alberta Innovates and the Forest
401 Resources Improvement Association of Alberta. M-JF acknowledges funding from NSERC Discovery
402 Grant and NSERC CRC. RR is funded by New Zealand's Biological Heritage Ngā Koiora Tuku Iho
403 National Science Challenge, administered by New Zealand Ministry of Business, Innovation, and
404 Employment. BM is funded by the NSERC Alexander Graham Bell Canada Graduate Scholarship and the
405 FRQNT master's scholarship. LP acknowledges funding from NSERC Discovery Grant (NSERC
406 RGPIN-2019-05771). TP acknowledges financial support from NSERC through the Discovery Grants and
407 Discovery Accelerator Supplement programs.

408 **Conflict of interest:** The authors have no conflict interests to disclose

409 **Authors' contributions:** TS, SB, and TP designed the study and performed the analysis; GVDR, MF, and
410 RR provided additional feedback on the analyses. DC, BM, and FB helped with data collection. All
411 authors contributed to writing and editing the manuscript.

412 **Data availability:** All code and data used in this manuscript is publicly available and archived on OSF
413 <https://osf.io/2zwqm/> and is currently referenced in the manuscript.

414 References

- 415 Adam, R. (2014). *Elephant treaties: The Colonial legacy of the biodiversity crisis*. UPNE.
- 416 Albouy, C., Archambault, P., Appeltans, W., Araújo, M. B., Beauchesne, D., Cazelles, K., Cirtwill, A. R.,
417 Fortin, M.-J., Galiana, N., Leroux, S. J., Pellissier, L., Poisot, T., Stouffer, D. B., Wood, S. A., & Gravel, D.
418 (2019). The marine fish food web is globally connected. *Nature Ecology & Evolution*, 3(8, 8),
419 1153–1161. <https://doi.org/10.1038/s41559-019-0950-y>
- 420 Arsov, N., & Mirceva, G. (2019, November 26). *Network Embedding: An Overview*.
421 <http://arxiv.org/abs/1911.11726>
- 422 Banville, F., Vissault, S., & Poisot, T. (2021). Mangal.jl and EcologicalNetworks.jl: Two complementary
423 packages for analyzing ecological networks in Julia. *Journal of Open Source Software*, 6(61), 2721.
424 <https://doi.org/10.21105/joss.02721>

- 425 Bascompte, J., Jordano, P., Melian, C. J., & Olesen, J. M. (2003). The nested assembly of plant-animal
426 mutualistic networks. *Proceedings of the National Academy of Sciences*, 100(16), 9383–9387.
427 <https://doi.org/10.1073/pnas.1633576100>
- 428 Beckerman, A. P., Petchey, O. L., & Warren, P. H. (2006). Foraging biology predicts food web complexity.
429 *Proceedings of the National Academy of Sciences*, 103(37), 13745–13749.
430 <https://doi.org/10.1073/pnas.0603039103>
- 431 Bezanson, J., Edelman, A., Karpinski, S., & Shah, V. (2017). Julia: A Fresh Approach to Numerical
432 Computing. *SIAM Review*, 59(1), 65–98. <https://doi.org/10.1137/141000671>
- 433 Boeckaerts, D., Stock, M., Criel, B., Gerstmans, H., De Baets, B., & Briers, Y. (2021). Predicting
434 bacteriophage hosts based on sequences of annotated receptor-binding proteins. *Scientific Reports*,
435 11(1, 1), 1467. <https://doi.org/10.1038/s41598-021-81063-4>
- 436 Braga, M. P., Janz, N., Nylin, S., Ronquist, F., & Landis, M. J. (2021). Phylogenetic reconstruction of
437 ancestral ecological networks through time for pierid butterflies and their host plants. *Ecology Letters*,
438 n/a(n/a). <https://doi.org/10.1111/ele.13842>
- 439 Brose, U., Jonsson, T., Berlow, E. L., Warren, P., Banasek-Richter, C., Bersier, L.-F., Blanchard, J. L., Brey,
440 T., Carpenter, S. R., Blandenier, M.-F. C., Cushing, L., Dawah, H. A., Dell, T., Edwards, F.,
441 Harper-Smith, S., Jacob, U., Ledger, M. E., Martinez, N. D., Memmott, J., ... Cohen, J. E. (2006).
442 ConsumerResource Body-Size Relationships in Natural Food Webs. *Ecology*, 87(10), 2411–2417.
443 [https://doi.org/10.1890/0012-9658\(2006\)87%5B2411:CBRINF%5D2.0.CO;2](https://doi.org/10.1890/0012-9658(2006)87%5B2411:CBRINF%5D2.0.CO;2)
- 444 Buxton, R. T., Bennett, J. R., Reid, A. J., Shulman, C., Cooke, S. J., Francis, C. M., Nyboer, E. A., Pritchard,
445 G., Binley, A. D., Avery-Gomm, S., Ban, N. C., Beazley, K. F., Bennett, E., Blight, L. K., Bortolotti, L. E.,
446 Camfield, A. F., Gadallah, F., Jacob, A. L., Naujokaitis-Lewis, I., ... Smith, P. A. (2021). Key
447 information needs to move from knowledge to action for biodiversity conservation in Canada.
448 *Biological Conservation*, 256, 108983. <https://doi.org/10.1016/j.biocon.2021.108983>
- 449 Cai, H., Zheng, V. W., & Chang, K. C.-C. (2017). *A Comprehensive Survey of Graph Embedding: Problems,*
450 *Techniques and Applications*. <http://arxiv.org/abs/1709.07604>
- 451 Cameron, E. K., Sundqvist, M. K., Keith, S. A., CaraDonna, P. J., Mousing, E. A., Nilsson, K. A., Metcalfe,
452 D. B., & Classen, A. T. (2019). Uneven global distribution of food web studies under climate change.
453 *Ecosphere*, 10(3), e02645. <https://doi.org/10.1002/ecs2.2645>

- 454 Cao, R.-M., Liu, S.-Y., & Xu, X.-K. (2019). Network embedding for link prediction: The pitfall and
455 improvement. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 29(10), 103102.
456 <https://doi.org/10.1063/1.5120724>
- 457 Cavender-Bares, J., Kozak, K. H., Fine, P. V. A., & Kembel, S. W. (2009). The merging of community
458 ecology and phylogenetic biology. *Ecology Letters*, 12(7), 693–715.
459 <https://doi.org/10.1111/j.1461-0248.2009.01314.x>
- 460 Choudry, A. (2013). Saving biodiversity, for whom and for what? Conservation NGOs, complicity,
461 colonialism and conquest in an era of capitalist globalization. In *NGOization: Complicity,
462 contradictions and prospects* (pp. 24–44). Bloomsbury Publishing.
- 463 Cirtwill, A. R., Ekl, A., Roslin, T., Wootton, K., & Gravel, D. (2019). A quantitative framework for
464 investigating the reliability of empirical network construction. *Methods in Ecology and Evolution*, 0.
465 <https://doi.org/10.1111/2041-210X.13180>
- 466 Cirtwill, A. R., & Hambäck, P. (2021). Building food networks from molecular data: Bayesian or
467 fixed-number thresholds for including links. *Basic and Applied Ecology*, 50, 67–76.
468 <https://doi.org/10.1016/j.baae.2020.11.007>
- 469 Csermely, P. (2004). Strong links are important, but weak links stabilize them. *Trends in Biochemical
470 Sciences*, 29(7), 331–334. <https://doi.org/10.1016/j.tibs.2004.05.004>
- 471 Dalla Riva, G. V., & Stouffer, D. B. (2016). Exploring the evolutionary signature of food webs' backbones
472 using functional traits. *Oikos*, 125(4), 446–456. <https://doi.org/10.1111/oik.02305>
- 473 Dallas, T., Park, A. W., & Drake, J. M. (2017). Predicting cryptic links in host-parasite networks. *PLOS
474 Computational Biology*, 13(5), e1005557. <https://doi.org/10.1371/journal.pcbi.1005557>
- 475 Dansereau, G., & Poisot, T. (2021). SimpleSDMLayers.jl and GBIF.jl: A Framework for Species
476 Distribution Modeling in Julia. *Journal of Open Source Software*, 6(57), 2872.
477 <https://doi.org/10.21105/joss.02872>
- 478 Domínguez, L., & Luoma, C. (2020). Decolonising Conservation Policy: How Colonial Land and
479 Conservation Ideologies Persist and Perpetuate Indigenous Injustices at the Expense of the
480 Environment. *Land*, 9(3, 3), 65. <https://doi.org/10.3390/land9030065>
- 481 Dormann, C. F., Gruber, B., Winter, M., & Herrmann, D. (2010). Evolution of climate niches in European

- 482 mammals? *Biology Letters*, 6(2), 229–232. <https://doi.org/10.1098/rsbl.2009.0688>
- 483 Dunne, J. A. (2006). The Network Structure of Food Webs. In J. A. Dunne & M. Pascual (Eds.), *Ecological*
484 *networks: Linking structure and dynamics* (pp. 27–86). Oxford University Press.
- 485 Eichhorn, M. P., Baker, K., & Griffiths, M. (2019). Steps towards decolonising biogeography. *Frontiers of*
486 *Biogeography*, 12(1), 1–7. <https://doi.org/10.21425/F5FBG44795>
- 487 Eklöf, A., & Stouffer, D. B. (2016). The phylogenetic component of food web structure and intervality.
488 *Theoretical Ecology*, 9(1), 107–115. <https://doi.org/10.1007/s12080-015-0273-9>
- 489 Fortin, M.-J., Dale, M. R. T., & Brimacombe, C. (2021). Network ecology in dynamic landscapes.
490 *Proceedings of the Royal Society B: Biological Sciences*, 288(1949), rspb.2020.1889, 20201889.
491 <https://doi.org/10.1098/rspb.2020.1889>
- 492 Galiana, N., Barros, C., Braga, J., Ficetola, G. F., Maiorano, L., Thuiller, W., Montoya, J. M., & Lurgi, M.
493 (2021). The spatial scaling of food web structure across European biogeographical regions. *Ecography*,
494 n/a(n/a). <https://doi.org/10.1111/ecog.05229>
- 495 Galiana, N., Hawkins, B. A., & Montoya, J. M. (2019). The geographical variation of network structure is
496 scale dependent: Understanding the biotic specialization of hostparasitoid networks. *Ecography*, 42(6),
497 1175–1187. <https://doi.org/10.1111/ecog.03684>
- 498 Galiana, N., Lurgi, M., Claramunt-López, B., Fortin, M.-J., Leroux, S., Cazelles, K., Gravel, D., & Montoya,
499 J. M. (2018). The spatial scaling of species interaction networks. *Nature Ecology & Evolution*, 2(5),
500 782–790. <https://doi.org/10.1038/s41559-018-0517-3>
- 501 Garland, T., JR., Midford, P. E., & Ives, A. R. (1999). An Introduction to Phylogenetically Based Statistical
502 Methods, with a New Method for Confidence Intervals on Ancestral Values1. *American Zoologist*,
503 39(2), 374–388. <https://doi.org/10.1093/icb/39.2.374>
- 504 Garlaschelli, D., Hollander, F. den, & Roccaverde, A. (2018). Covariance structure behind breaking of
505 ensemble equivalence in random graphs. *Journal of Statistical Physics*, 173(3-4), 644–662.
506 <https://doi.org/10.1007/s10955-018-2114-x>
- 507 GBIF Secretariat. (2021). *GBIF Backbone Taxonomy*. <https://doi.org/10.15468/39omei>
- 508 Gerhold, P., Cahill, J. F., Winter, M., Bartish, I. V., & Prinzing, A. (2015). Phylogenetic patterns are not
509 proxies of community assembly mechanisms (they are far better). *Functional Ecology*, 29(5), 600–614.

- 510 <https://doi.org/10.1111/1365-2435.12425>
- 511 Gravel, D., Baiser, B., Dunne, J. A., Kopalke, J.-P., Martinez, N. D., Nyman, T., Poisot, T., Stouffer, D. B.,
512 Tylianakis, J. M., Wood, S. A., & Roslin, T. (2018). Bringing Elton and Grinnell together: A quantitative
513 framework to represent the biogeography of ecological interaction networks. *Ecography*, 0(0).
- 514 <https://doi.org/10.1111/ecog.04006>
- 515 Grünig, M., Mazzi, D., Calanca, P., Karger, D. N., & Pellissier, L. (2020). Crop and forest pest metawebs
516 shift towards increased linkage and suitability overlap under climate change. *Communications Biology*,
517 3(1, 1), 1–10. <https://doi.org/10.1038/s42003-020-0962-9>
- 518 Halevy, A., Norvig, P., & Pereira, F. (2009). The Unreasonable Effectiveness of Data. *IEEE Intelligent
519 Systems*, 24(2), 8–12. <https://doi.org/10.1109/MIS.2009.36>
- 520 Halko, N., Martinsson, P. G., & Tropp, J. A. (2011). Finding Structure with Randomness: Probabilistic
521 Algorithms for Constructing Approximate Matrix Decompositions. *SIAM Review*, 53(2), 217–288.
522 <https://doi.org/10.1137/090771806>
- 523 Herbert, F. (1965). *Dune* (1st ed.). Chilton Book Company.
- 524 Hoffmann, J., Bar-Sinai, Y., Lee, L. M., Andrejevic, J., Mishra, S., Rubinstein, S. M., & Rycroft, C. H. (2019).
525 Machine learning in a data-limited regime: Augmenting experiments with synthetic data uncovers
526 order in crumpled sheets. *Science Advances*, 5(4), eaau6792.
527 <https://doi.org/10.1126/sciadv.aau6792>
- 528 Holm, E. A. (2019). In defense of the black box. *Science*, 364(6435), 26–27.
529 <https://doi.org/10.1126/science.aax0162>
- 530 Hortal, J., de Bello, F., Diniz-Filho, J. A. F., Lewinsohn, T. M., Lobo, J. M., & Ladle, R. J. (2015). Seven
531 Shortfalls that Beset Large-Scale Knowledge of Biodiversity. *Annual Review of Ecology, Evolution, and
532 Systematics*, 46(1), 523–549. <https://doi.org/10.1146/annurev-ecolsys-112414-054400>
- 533 Hutchinson, M. C., Cagua, E. F., & Stouffer, D. B. (2017). Cophylogenetic signal is detectable in pollination
534 interactions across ecological scales. *Ecology*, n/a–n/a. <https://doi.org/10.1002/ecy.1955>
- 535 Jordano, P. (2016a). Chasing Ecological Interactions. *PLOS Biol*, 14(9), e1002559.
536 <https://doi.org/10.1371/journal.pbio.1002559>

- 537 Jordano, P. (2016b). Sampling networks of ecological interactions. *Functional Ecology*, 30(12), 1883–1893.
- 538 <https://doi.org/10.1111/1365-2435.12763>
- 539 Kawatsu, K., Ushio, M., van Veen, F. J. F., & Kondoh, M. (2021). Are networks of trophic interactions
540 sufficient for understanding the dynamics of multi-trophic communities? Analysis of a tri-trophic
541 insect food-web time-series. *Ecology Letters*, 24(3), 543–552. <https://doi.org/10.1111/ele.13672>
- 542 Kéfi, S., Berlow, E. L., Wieters, E. A., Navarrete, S. A., Petchey, O. L., Wood, S. A., Boit, A., Joppa, L. N.,
543 Lafferty, K. D., Williams, R. J., Martinez, N. D., Menge, B. A., Blanchette, C. A., Iles, A. C., & Brose, U.
544 (2012). More than a meal... integrating non-feeding interactions into food webs: More than a meal
545 *Ecology Letters*, 15(4), 291–300. <https://doi.org/10.1111/j.1461-0248.2011.01732.x>
- 546 Lamba, A., Cassey, P., Segaran, R. R., & Koh, L. P. (2019). Deep learning for environmental conservation.
547 *Current Biology*, 29(19), R977–R982. <https://doi.org/10.1016/j.cub.2019.08.016>
- 548 Litsios, G., & Salamin, N. (2012). Effects of Phylogenetic Signal on Ancestral State Reconstruction.
549 *Systematic Biology*, 61(3), 533–538. <https://doi.org/10.1093/sysbio/syr124>
- 550 Machen, R., & Nost, E. (2021). Thinking algorithmically: The making of hegemonic knowledge in climate
551 governance. *Transactions of the Institute of British Geographers*, 46(3), 555–569.
552 <https://doi.org/10.1111/tran.12441>
- 553 Maiorano, L., Montemaggiori, A., Ficetola, G. F., O'Connor, L., & Thuiller, W. (2020a). TETRA-EU 1.0: A
554 species-level trophic metaweb of European tetrapods. *Global Ecology and Biogeography*, 29(9),
555 1452–1457. <https://doi.org/10.1111/geb.13138>
- 556 Maiorano, L., Montemaggiori, A., Ficetola, G. F., O'Connor, L., & Thuiller, W. (2020b). *Data from:
557 Tetra-EU 1.0: A species-level trophic meta-web of European tetrapods* (Version 3, pp. 16596876 bytes)
558 [Data set]. Dryad. <https://doi.org/10.5061/DRYAD.JM63XSJ7B>
- 559 Marco, M. D., Venter, O., Possingham, H. P., & Watson, J. E. M. (2018). Changes in human footprint drive
560 changes in species extinction risk. *Nature Communications*, 9(1), 4621.
561 <https://doi.org/10.1038/s41467-018-07049-5>
- 562 McLeod, A., Leroux, S. J., Gravel, D., Chu, C., Cirtwill, A. R., Fortin, M.-J., Galiana, N., Poisot, T., & Wood,
563 S. A. (2021). Sampling and asymptotic network properties of spatial multi-trophic networks. *Oikos*,
564 n/a(n/a). <https://doi.org/10.1111/oik.08650>

- 565 Mora, B. B., Gravel, D., Gilarranz, L. J., Poisot, T., & Stouffer, D. B. (2018). Identifying a common backbone
566 of interactions underlying food webs from different ecosystems. *Nature Communications*, 9(1), 2603.
567
<https://doi.org/10.1038/s41467-018-05056-0>
- 568 Morales-Castilla, I., Matias, M. G., Gravel, D., & Araújo, M. B. (2015). Inferring biotic interactions from
569 proxies. *Trends in Ecology & Evolution*, 30(6), 347–356.
570
<https://doi.org/10.1016/j.tree.2015.03.014>
- 571 Mosebo Fernandes, A. C., Quintero Gonzalez, R., Lenihan-Clarke, M. A., Leslie Trotter, E. F., & Jokar
572 Arsanjani, J. (2020). Machine Learning for Conservation Planning in a Changing Climate.
573 *Sustainability*, 12(18, 18), 7657. <https://doi.org/10.3390/su12187657>
- 574 Mouquet, N., Devictor, V., Meynard, C. N., Munoz, F., Bersier, L.-F., Chave, J., Couteron, P., Dalecky, A.,
575 Fontaine, C., Gravel, D., Hardy, O. J., Jabot, F., Lavergne, S., Leibold, M., Mouillot, D., Münkemüller,
576 T., Pavoine, S., Prinzing, A., Rodrigues, A. S. L., ... Thuiller, W. (2012). Ecophylogenetics: Advances
577 and perspectives. *Biological Reviews*, 87(4), 769–785.
578
<https://doi.org/10.1111/j.1469-185X.2012.00224.x>
- 579 Nenzén, H. K., Montoya, D., & Varela, S. (2014). The Impact of 850,000 Years of Climate Changes on the
580 Structure and Dynamics of Mammal Food Webs. *PLOS ONE*, 9(9), e106651.
581
<https://doi.org/10.1371/journal.pone.0106651>
- 582 Neutel, A.-M., Heesterbeek, J. A. P., & de Ruiter, P. C. (2002). Stability in Real Food Webs: Weak Links in
583 Long Loops. *Science*, 296(5570), 1120–1123. <https://doi.org/10.1126/science.1068326>
- 584 No'kmaq, M., Marshall, A., Beazley, K. F., Hum, J., joudry, shalan, Papadopoulos, A., Pictou, S., Rabesca,
585 J., Young, L., & Zurba, M. (2021). “Awakening the sleeping giant”: Re-Indigenization principles for
586 transforming biodiversity conservation in Canada and beyond. *FACETS*, 6(1), 839–869.
- 587 Nost, E., & Goldstein, J. E. (2021). A political ecology of data. *Environment and Planning E: Nature and
588 Space*, 25148486211043503. <https://doi.org/10.1177/25148486211043503>
- 589 O'Connor, L. M. J., Pollock, L. J., Braga, J., Ficetola, G. F., Maiorano, L., Martinez-Almoyna, C.,
590 Montemaggioli, A., Ohlmann, M., & Thuiller, W. (2020). Unveiling the food webs of tetrapods across
591 Europe through the prism of the Eltonian niche. *Journal of Biogeography*, 47(1), 181–192.
592
<https://doi.org/10.1111/jbi.13773>

- 593 Pan, S. J., & Yang, Q. (2010). A Survey on Transfer Learning. *IEEE Transactions on Knowledge and Data
594 Engineering*, 22(10), 1345–1359. <https://doi.org/10.1109/TKDE.2009.191>
- 595 Park, J., & Newman, M. E. J. (2004). Statistical mechanics of networks. *Physical Review E*, 70(6), 066117.
596 <https://doi.org/10.1103/PhysRevE.70.066117>
- 597 Pedersen, E. J., Thompson, P. L., Ball, R. A., Fortin, M.-J., Gouhier, T. C., Link, H., Moritz, C., Nenzen, H.,
598 Stanley, R. R. E., Taranu, Z. E., Gonzalez, A., Guichard, F., & Pepin, P. (2017). Signatures of the
599 collapse and incipient recovery of an overexploited marine ecosystem. *Royal Society Open Science*, 4(7),
600 170215. <https://doi.org/10.1098/rsos.170215>
- 601 Perretti, C. T., Munch, S. B., & Sugihara, G. (2013). Model-free forecasting outperforms the correct
602 mechanistic model for simulated and experimental data. *Proceedings of the National Academy of
603 Sciences*, 110(13), 5253–5257. <https://doi.org/10.1073/pnas.1216076110>
- 604 Pires, M. M., Koch, P. L., Fariña, R. A., de Aguiar, M. A. M., dos Reis, S. F., & Guimarães, P. R. (2015).
605 Pleistocene megafaunal interaction networks became more vulnerable after human arrival.
606 *Proceedings of the Royal Society B: Biological Sciences*, 282(1814), 20151367.
607 <https://doi.org/10.1098/rspb.2015.1367>
- 608 Poelen, J. H., Simons, J. D., & Mungall, C. J. (2014). Global biotic interactions: An open infrastructure to
609 share and analyze species-interaction datasets. *Ecological Informatics*, 24, 148–159.
610 <https://doi.org/10.1016/j.ecoinf.2014.08.005>
- 611 Poisot, T., Belisle, Z., Hoebelke, L., Stock, M., & Szefer, P. (2019). EcologicalNetworks.jl - analysing
612 ecological networks. *Ecography*. <https://doi.org/10.1111/ecog.04310>
- 613 Poisot, T., Bergeron, G., Cazelles, K., Dallas, T., Gravel, D., MacDonald, A., Mercier, B., Violet, C., &
614 Vissault, S. (2021). Global knowledge gaps in species interaction networks data. *Journal of
615 Biogeography*, n/a(n/a). <https://doi.org/10.1111/jbi.14127>
- 616 Poisot, T., Cirtwill, A. R., Cazelles, K., Gravel, D., Fortin, M.-J., & Stouffer, D. B. (2016). The structure of
617 probabilistic networks. *Methods in Ecology and Evolution*, 7(3), 303–312.
618 <https://doi.org/10.1111/2041-210X.12468>
- 619 Poisot, T., Ouellet, M.-A., Mollentze, N., Farrell, M. J., Becker, D. J., Albery, G. F., Gibb, R. J., Seifert, S. N.,
620 & Carlson, C. J. (2021, May 31). *Imputing the mammalian virome with linear filtering and singular*

- 621 value decomposition. <http://arxiv.org/abs/2105.14973>
- 622 Poisot, T., & Stouffer, D. B. (2018). Interactions retain the co-phylogenetic matching that communities lost.
- 623 *Oikos*, 127(2), 230–238. <https://doi.org/10.1111/oik.03788>
- 624 Poisot, T., Stouffer, D. B., & Gravel, D. (2015). Beyond species: Why ecological interaction networks vary
- 625 through space and time. *Oikos*, 124(3), 243–251. <https://doi.org/10.1111/oik.01719>
- 626 Price, P. W. (2003). *Macroevolutionary theory on macroecological patterns*. Cambridge University Press.
- 627 Ray, J. C., Grimm, J., & Olive, A. (2021). The biodiversity crisis in Canada: Failures and challenges of
- 628 federal and sub-national strategic and legal frameworks. *FACETS*, 6, 1044–1068.
- 629 <https://doi.org/10.1139/facets-2020-0075>
- 630 Reeve, R., Leinster, T., Cobbold, C. A., Thompson, J., Brummitt, N., Mitchell, S. N., & Matthews, L. (2016,
- 631 December 8). *How to partition diversity*. <http://arxiv.org/abs/1404.6520>
- 632 Rosado, B. H. P., Dias, A., & de Mattos, E. (2013). Going Back to Basics: Importance of Ecophysiology
- 633 when Choosing Functional Traits for Studying Communities and Ecosystems. *Natureza &*
- 634 *Conservação Revista Brasileira de Conservação Da Natureza*, 11, 15–22.
- 635 <https://doi.org/10.4322/natcon.2013.002>
- 636 Runghen, R., Stouffer, D. B., & Dalla Riva, G. V. (2021). *Exploiting node metadata to predict interactions in*
- 637 *large networks using graph embedding and neural networks*.
- 638 <https://doi.org/10.1101/2021.06.10.447991>
- 639 Schoener, T. W. (1989). Food webs from the small to the large. *Ecology*, 70(6), 1559–1589.
- 640 Solís-Lemus, C., Bastide, P., & Ané, C. (2017). PhyloNetworks: A Package for Phylogenetic Networks.
- 641 *Molecular Biology and Evolution*, 34(12), 3292–3298. <https://doi.org/10.1093/molbev/msx235>
- 642 Stock, M. (2021). Pairwise learning for predicting pollination interactions based on traits and phylogeny.
- 643 *Ecological Modelling*, 14.
- 644 Stouffer, D. B., Sales-Pardo, M., Sirer, M. I., & Bascompte, J. (2012). Evolutionary Conservation of Species'
- 645 Roles in Food Webs. *Science*, 335(6075), 1489–1492. <https://doi.org/10.1126/science.1216556>
- 646 Strong, J. S., & Leroux, S. J. (2014). Impact of Non-Native Terrestrial Mammals on the Structure of the
- 647 Terrestrial Mammal Food Web of Newfoundland, Canada. *PLOS ONE*, 9(8), e106264.
- 648 <https://doi.org/10.1371/journal.pone.0106264>

- 649 Strydom, T., Catchen, M. D., Banville, F., Caron, D., Dansereau, G., Desjardins-Proulx, P., Forero-Muñoz,
650 N. R., Higino, G., Mercier, B., Gonzalez, A., Gravel, D., Pollock, L., & Poisot, T. (2021). A roadmap
651 towards predicting species interaction networks (across space and time). *Philosophical Transactions of*
652 *the Royal Society B: Biological Sciences*, 376(1837), 20210063.
653 <https://doi.org/10.1098/rstb.2021.0063>
- 654 Strydom, T., Dalla Riva, G. V., & Poisot, T. (2021). SVD Entropy Reveals the High Complexity of Ecological
655 Networks. *Frontiers in Ecology and Evolution*, 9. <https://doi.org/10.3389/fevo.2021.623141>
- 656 Thessen, A. E., & Parr, C. S. (2014). Knowledge extraction and semantic annotation of text from the
657 encyclopedia of life. *PloS One*, 9(3), e89550.
- 658 Torrey, L., & Shavlik, J. (2010). Transfer learning. In *Handbook of research on machine learning*
659 *applications and trends: Algorithms, methods, and techniques* (pp. 242–264). IGI global.
- 660 Trøjelsgaard, K., & Olesen, J. M. (2016). Ecological networks in motion: Micro- and macroscopic
661 variability across scales. *Functional Ecology*, 30(12), 1926–1935.
662 <https://doi.org/10.1111/1365-2435.12710>
- 663 Turak, E., Brazill-Boast, J., Cooney, T., Drielsma, M., DelaCruz, J., Dunkerley, G., Fernandez, M., Ferrier,
664 S., Gill, M., Jones, H., Koen, T., Leys, J., McGeoch, M., Mihoub, J.-B., Scanes, P., Schmeller, D., &
665 Williams, K. (2017). Using the essential biodiversity variables framework to measure biodiversity
666 change at national scale. *Biological Conservation*, 213, 264–271.
667 <https://doi.org/10.1016/j.biocon.2016.08.019>
- 668 Upham, N. S., Esselstyn, J. A., & Jetz, W. (2019). Inferring the mammal tree: Species-level sets of
669 phylogenies for questions in ecology, evolution, and conservation. *PLOS Biology*, 17(12), e3000494.
670 <https://doi.org/10.1371/journal.pbio.3000494>
- 671 van der Hoorn, P., Lippner, G., & Krioukov, D. (2018). Sparse Maximum-Entropy Random Graphs with a
672 Given Power-Law Degree Distribution. *Journal of Statistical Physics*, 173(3-4), 806–844.
673 <https://doi.org/10.1007/s10955-017-1887-7>
- 674 Vermaat, J. E., Dunne, J. A., & Gilbert, A. J. (2009). Major dimensions in food-web structure properties.
675 *Ecology*, 90(1), 278–282. <http://www.ncbi.nlm.nih.gov/pubmed/19294932>
- 676 Wood, S. A., Russell, R., Hanson, D., Williams, R. J., & Dunne, J. A. (2015). Effects of spatial scale of

- 677 sampling on food web structure. *Ecology and Evolution*, 5(17), 3769–3782.
- 678 <https://doi.org/10.1002/ece3.1640>
- 679 Xing, S., & Fayle, T. M. (2021). The rise of ecological network meta-analyses: Problems and prospects.
- 680 *Global Ecology and Conservation*, 30, e01805. <https://doi.org/10.1016/j.gecco.2021.e01805>
- 681 Yeakel, J. D., Pires, M. M., Rudolf, L., Dominy, N. J., Koch, P. L., Guimarães, P. R., & Gross, T. (2014).
- 682 Collapse of an ecological network in Ancient Egypt. *PNAS*, 111(40), 14472–14477.
- 683 <https://doi.org/10.1073/pnas.1408471111>
- 684 Youden, W. J. (1950). Index for rating diagnostic tests. *Cancer*, 3(1), 32–35.
- 685 [https://doi.org/10.1002/1097-0142\(1950\)3:1%3C32::AID-CNCR2820030106%3E3.0.CO;2-3](https://doi.org/10.1002/1097-0142(1950)3:1%3C32::AID-CNCR2820030106%3E3.0.CO;2-3)
- 686 Young, J.-G., Cantwell, G. T., & Newman, M. E. J. (2021). Bayesian inference of network structure from
- 687 unreliable data. *Journal of Complex Networks*, 8(6). <https://doi.org/10.1093/comnet/cnaa046>
- 688 Young, S. J., & Scheinerman, E. R. (2007). Random Dot Product Graph Models for Social Networks. In A.
- 689 Bonato & F. R. K. Chung (Eds.), *Algorithms and Models for the Web-Graph* (pp. 138–149). Springer.
- 690 https://doi.org/10.1007/978-3-540-77004-6_11
- 691 Zhu, M., & Ghodsi, A. (2006). Automatic dimensionality selection from the scree plot via the use of profile
- 692 likelihood. *Computational Statistics & Data Analysis*, 51(2), 918–930.
- 693 <https://doi.org/10.1016/j.csda.2005.09.010>

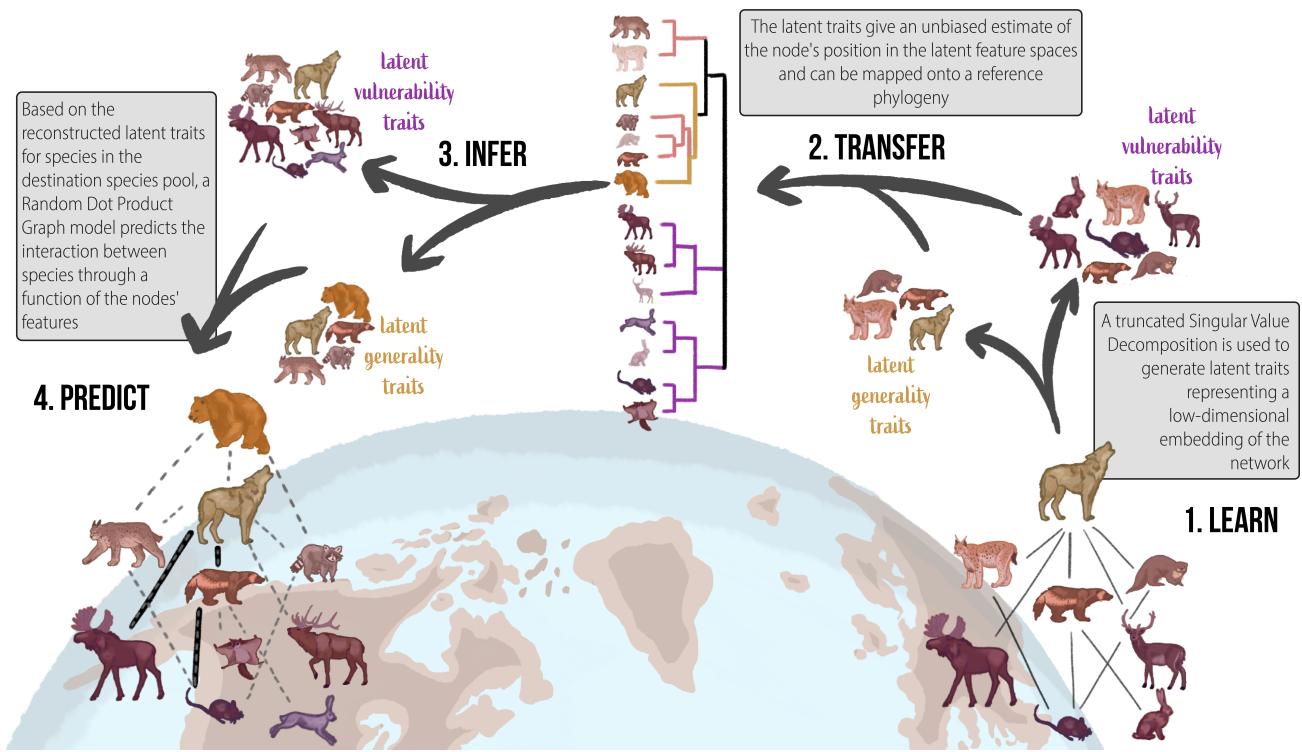


Figure 1: Overview of the phylogenetic transfer learning (and prediction) of species interaction networks. Starting from an initial, known, network, we learn its representation through a graph embedding step (here, a truncated Singular Value Decomposition; Step 1), yielding a series of latent traits (latent vulnerability traits are more representative of species at the lower trophic-level and latent generality traits are more representative of species at higher trophic-levels; *sensu* Schoener (1989)); second, for the destination species pool, we perform ancestral character estimation using a phylogeny (here, using a Brownian model for the latent traits; Step 2); we then sample from the reconstructed distribution of latent traits (Step 3) to generate a probabilistic metaweb at the destination (here, assuming a uniform distribution of traits), and threshold it to yield the final list of interactions (Step 4).

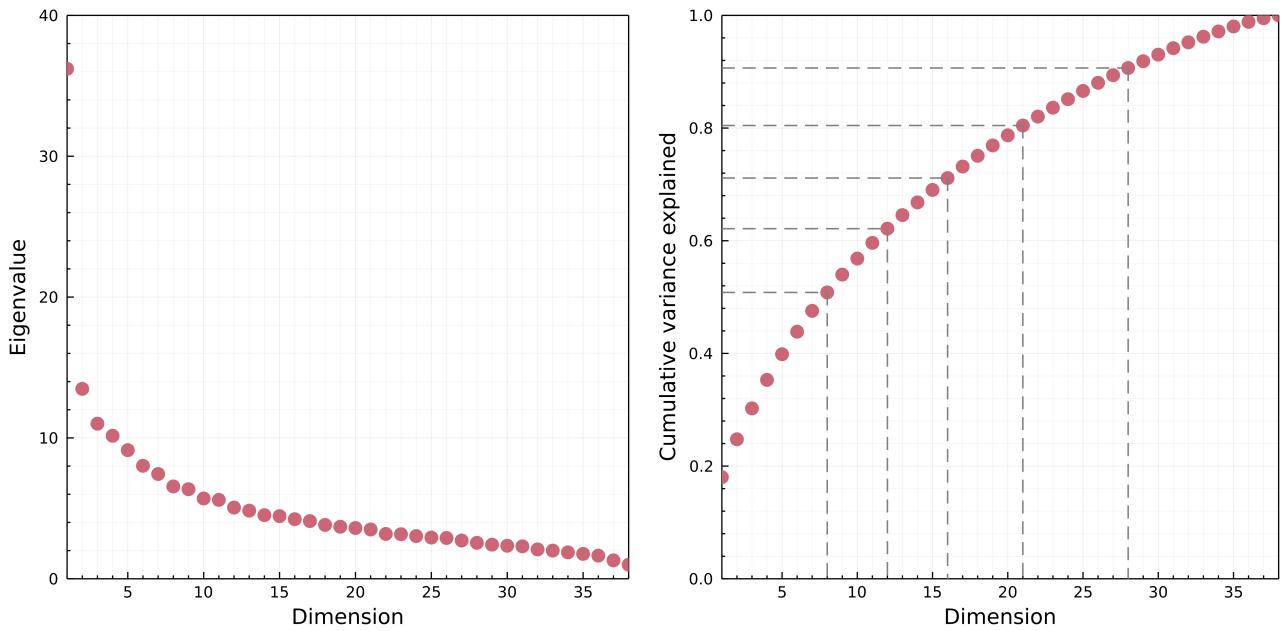


Figure 2: Left: representation of the scree plot of the singular values from the t-SVD on the European metaweb. The scree plot shows no obvious drop in the singular values that may be leveraged to automatically detect a minimal dimension for embedding, after e.g. Zhu & Ghodsi (2006). Right: cumulative fraction of variance explained by each dimension up to the rank of the European metaweb. The grey lines represent cutoffs at 50, 60, ..., 90% of variance explained. For the rest of the analysis, we reverted to an arbitrary threshold of 60% of variance explained, which represented a good tradeoff between accuracy and reduced number of features.

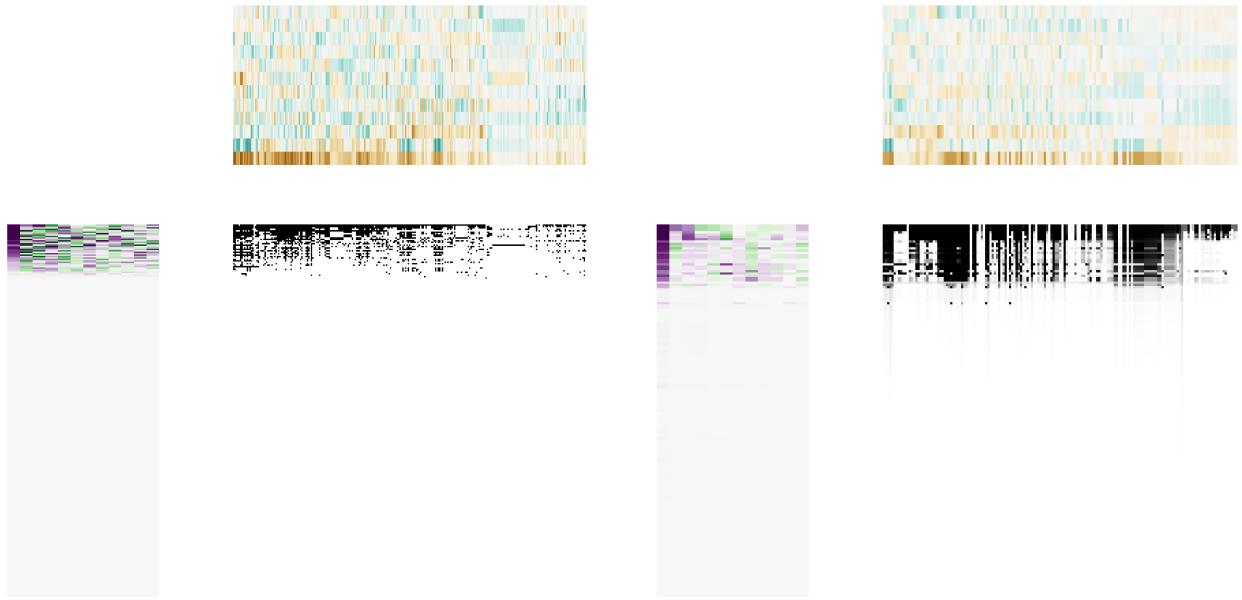


Figure 3: Visual representation of the left (green/purple) and right (green/brown) subspaces, alongside the adjacency matrix of the food web they encode (greyscale). The European metaweb is on the left, and the imputed Canadian metaweb (before data inflation) on the right. This figure illustrates how much structure the left subspace captures. As we show in fig. 6, the species with a value of 0 in the left subspace are species without any prey.

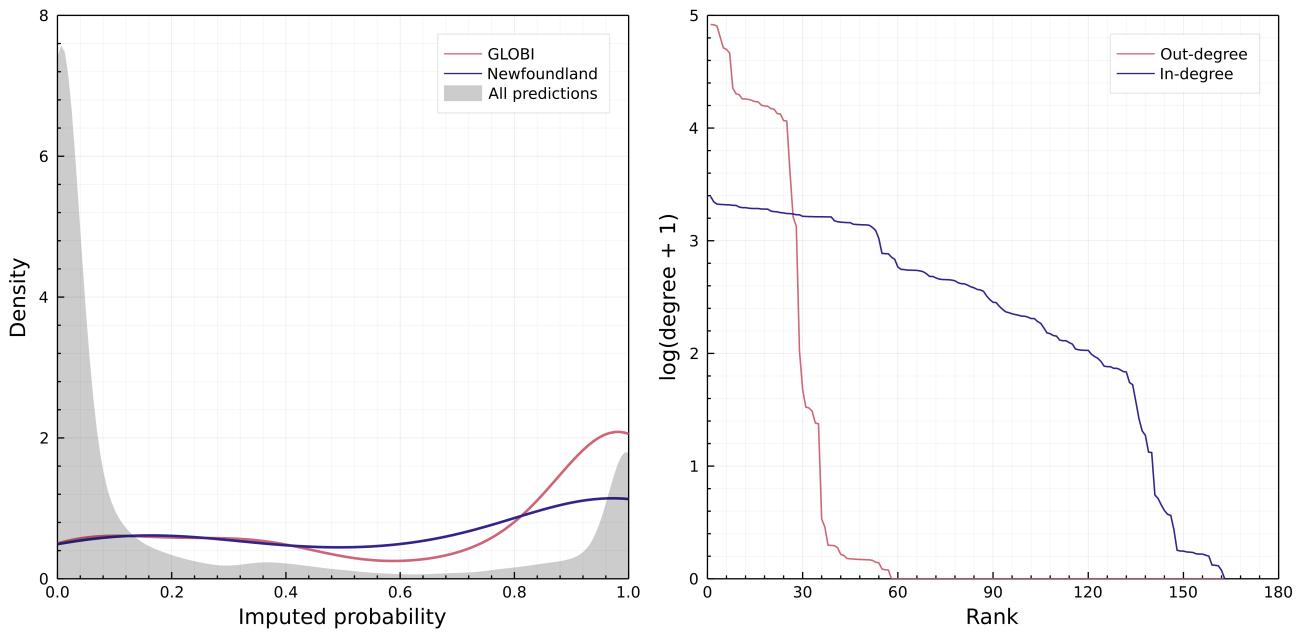


Figure 4: Left, comparison of the probabilities of interactions assigned by the model to all interactions (grey curve), the subset of interactions found in GLOBI (red), and in the Strong & Leroux (2014) Newfoundland dataset (blue). The model recovers more interactions with a low probability compared to data mining, which can suggest that collected datasets are biased towards more common or easy to identify interactions. Right, distribution of the in-degree and out-degree of the mammals from Canada in the reconstructed metaweb. This figure describes a flat, relatively short food web, in which there are few predators but a large number of preys.

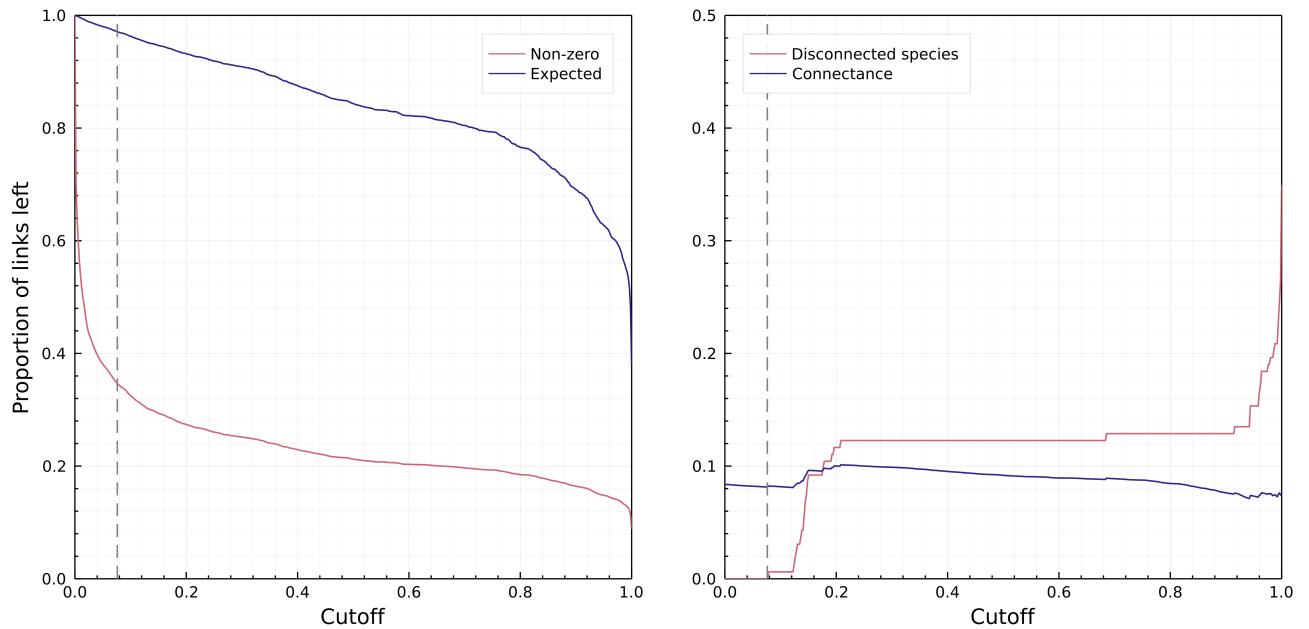


Figure 5: Left: effect of varying the cutoff for probabilities to be considered non-zero on the number of unique links and on \hat{L} , the probabilistic estimate of the number of links assuming that all interactions are independent. Right: effect of varying the cutoff on the number of disconnected species, and on network connectance. In both panels, the grey line indicates the cutoff $P(i \rightarrow j) \approx 0.08$ that resulted in the first species losing all of its interactions.

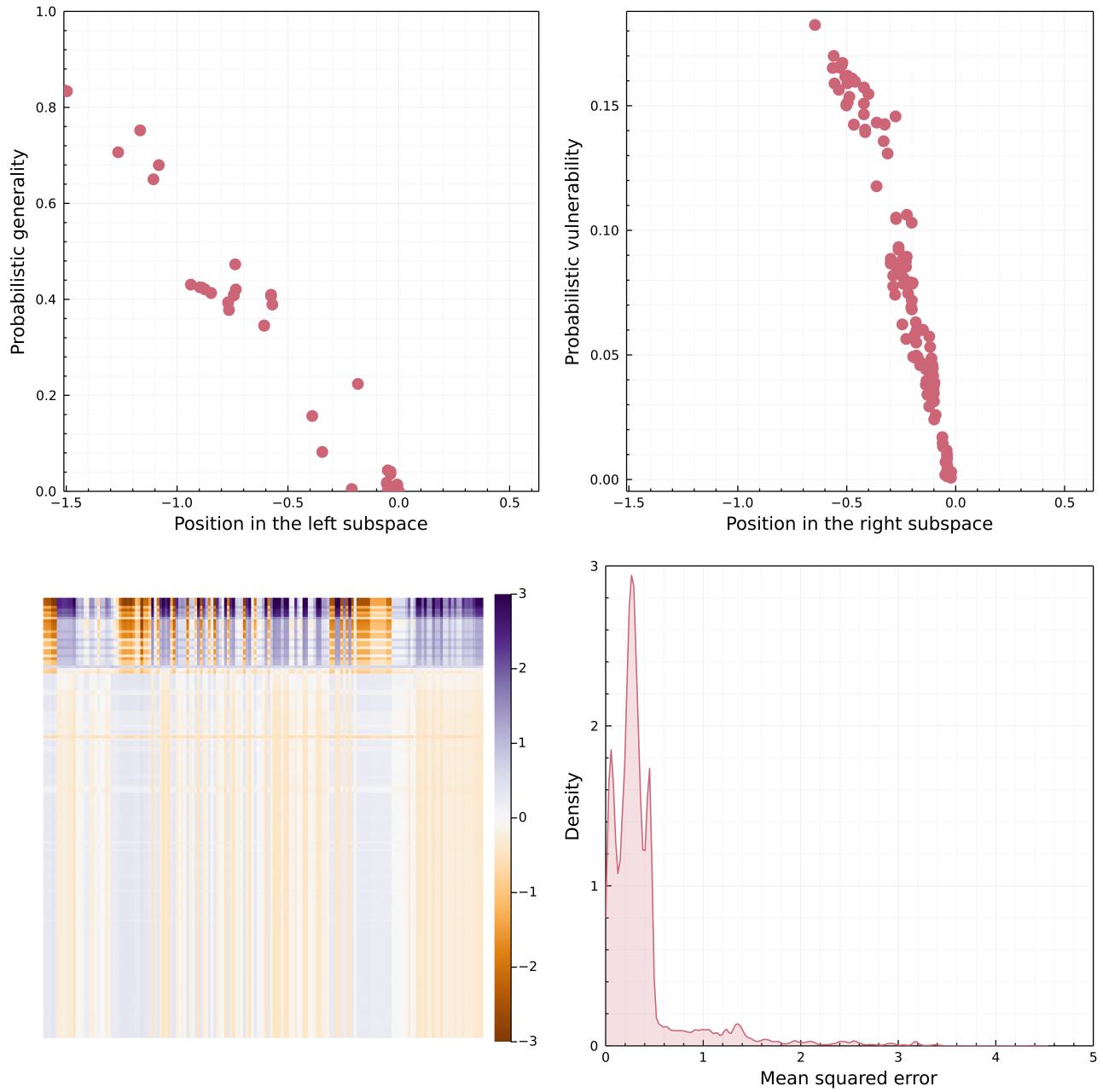


Figure 6: Top: biological significance of the first dimension. Left: there is a linear relationship between the values on the first dimension of the left subspace and the generality, *i.e.* the relative number of preys, *sensu* Schoener (1989). Species with a value of 0 in this subspace are at the bottom-most trophic level. Right: there is, similarly, a linear relationship between the position of a species on the first dimension of the right subspace and its vulnerability, *i.e.* the relative number of predators. Taken together, these two figures show that the first-order representation of this network would capture its degree distribution. Bottom: topological consequences of the first dimension. Left: differences in the z-score of the actual configuration model for the reconstructed network, and the prediction based only on the first dimension. Right: distribution of the differences in the left panel.